# Indexing for Near-Sorted Data

**Aneesh Raman**                    Subhadeep Sarkar
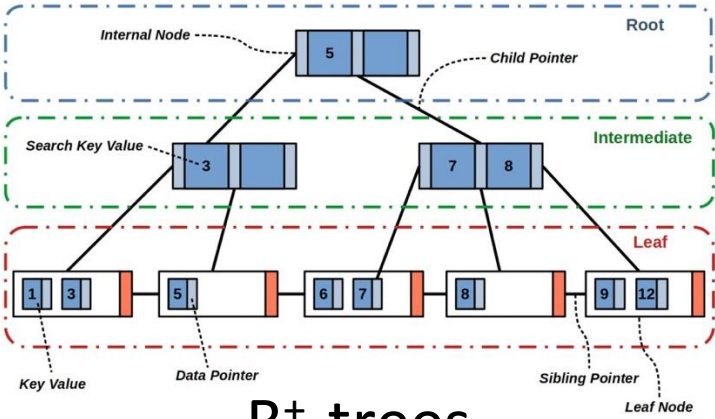
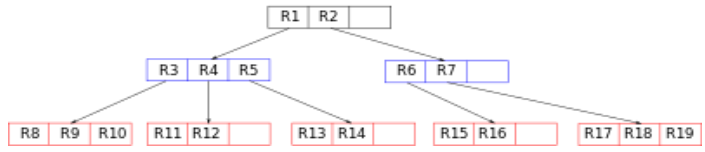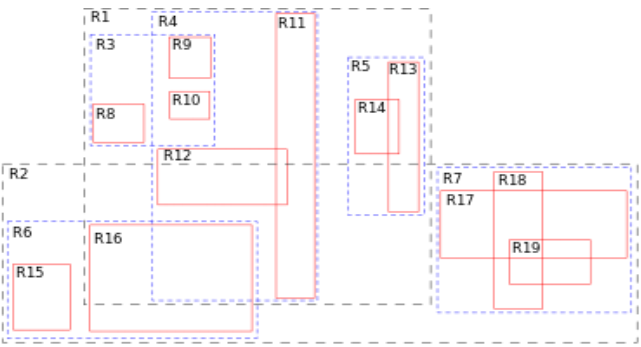Matthaios Olma                    Manos Athanassoulis
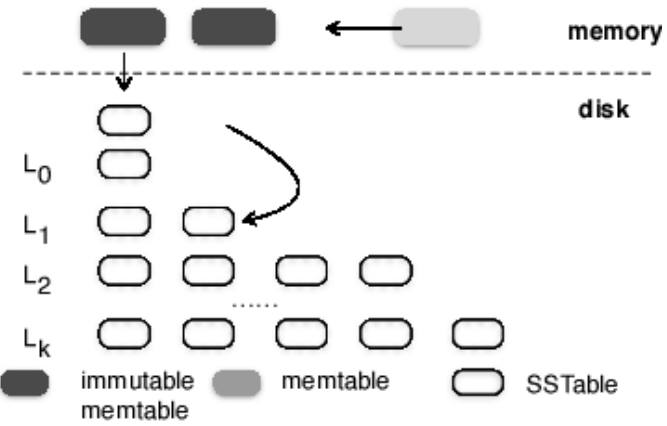
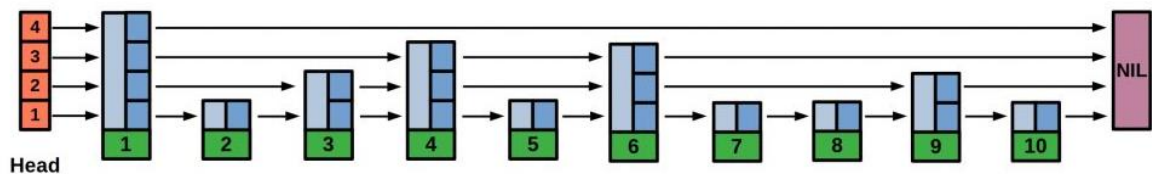BOSTON UNIVERSITY

DiSC lab

# Indexes in Databases



B⁺-trees



R-trees


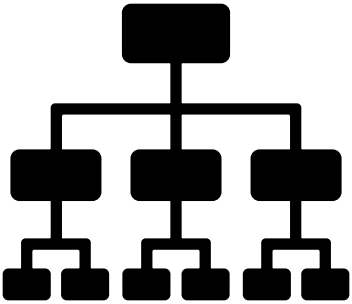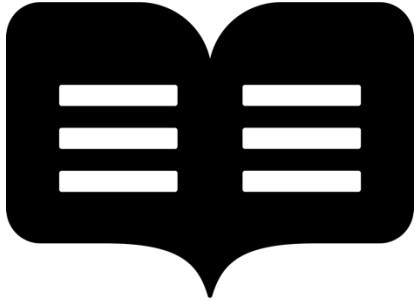
LSM-trees



Skip lists

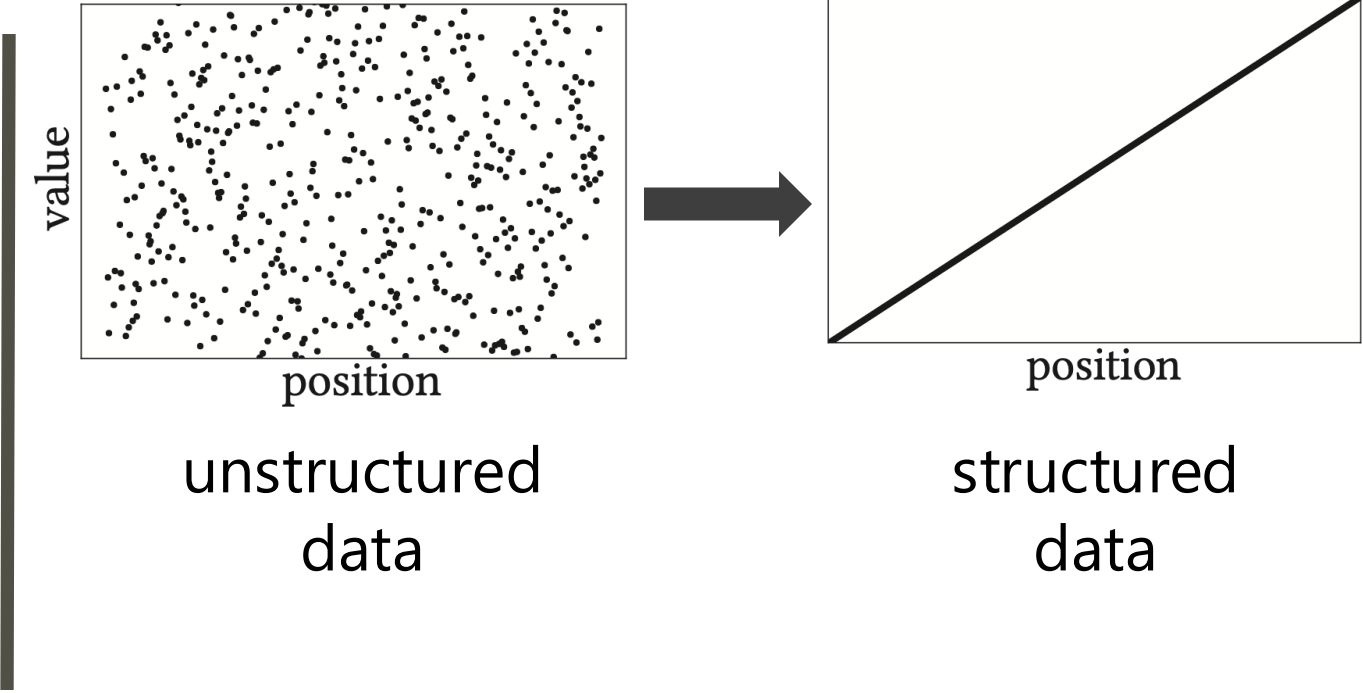# Indexes in Databases



organize
data

efficient
queries

unstructured
data

structured
data

**The process of inducing "*sortedness*" to an otherwise unsorted data collection**

What if data already has some structure?

What if data already has some structure?

Value / Position (time)

Value / Position (time)

≈ treated same as unstructured data!

Near-sorted data

# Irrespective of Sortedness, Same Ingestion Performance



Standard ingestion

Ingestion cost

Scrambled

Increasing data sortedness

Sorted

# Are There Faster Alternatives?

# Ideally, Higher Sortedness Should Lead to Faster Ingestion

# Near-Sorted Data is Frequently Found

Time Series

Stock market

Join/query

efficient reads + fast writes

classical indexes carry *redundant* **effort!**

# Vision for Sortedness-Aware Indexing

# Vision: Sortedness-Aware Indexes

# Vision: Sortedness-Aware Indexes

# Agenda

# Quantifying Data Sortedness

| Metric | Description |
|---|---|
| Inversions | # pairs in incorrect order |
| Runs | # increasing contiguous subsequences |
| Exchanges | least # swaps needed to establish total order |

# Any downsides of the "simple" metrics?

# Quantifying Data Sortedness

| Metric | Description |
|---|---|
| Inversions | # pairs in incorrect order |
| Runs | # increasing contiguous subsequences |
| Exchanges | least # swaps needed to establish total order |

| 6 | 7 | 8 | 9 | 10 | 1 | 2 | 3 | 4 | 5 |

# Quantifying Data Sortedness

| Metric | | Description |
|---|---|---|
| Inversions | ⊖ | # pairs in incorrect order |
| Runs | ✅ | # increasing contiguous subsequences |
| Exchanges | ⊖ | least # swaps needed to establish total order |

| 6 | 7 | 8 | 9 | 10 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|---|---|---|

global disorder

20

# Quantifying Data Sortedness

| Metric | Description |
|---|---|
| Inversions | # pairs in incorrect order |
| Runs | # increasing contiguous subsequences |
| Exchanges | least # swaps needed to establish total order |

| 2 | 1 | 4 | 3 | 6 | 5 | 8 | 7 | 10 | 9 |
|---|---|---|---|---|---|---|---|---|---|

# Quantifying Data Sortedness

| Metric | | Description |
|---|---|---|
| Inversions | ✅ | # pairs in incorrect order |
| Runs | ⛔ | # increasing contiguous subsequences |
| Exchanges | ✅ | least # swaps needed to establish total order |

| 2 | 1 | 4 | 3 | 6 | 5 | 8 | 7 | 10 | 9 |
|---|---|---|---|---|---|---|---|---|---|

local disorder

# (K, L)-Sortedness Metric

#. unordered entries    = K

| 1 | 8 | 3 | 4 | 5 | 6 | 7 | 2 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|----|

[inspired by BenMoshe, ICDT 2011]

# (K, L)-Sortedness Metric

#. unordered entries     = K

| 1 | 8 | 3 | 4 | 5 | 6 | 7 | 2 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|----|

max. displacement among unordered entries     = L

[inspired by BenMoshe, ICDT 2011]

# The Sortedness-Aware (SWARE) Paradigm

# Sortedness-Aware (SWARE) Paradigm



intelligent
buffering

+

opportunistic
bulk loading

+

increased fill
and split factor

**SWARE framework can be applied to any tree-index!**

# SWARE Ingestions

Buffer

Zonemap
(min-max)

700-749  750-799  800-849  850-910  870-950  920-991  890-1080  1090-1500

SWARE Buffer

flush 3 pages

tail leaf node

non-overlapping pages

non overlapping pages may move

flush non-overlapping pages to tree

B⁺-tree

Leaf pages  ...

BOSTON UNIVERSITY

DiSC lab

# SWARE Ingestions



Buffer

Zonemap (min-max)

700-749  750-799  800-849  850-910  870-950  920-991  890-1080  1090-1500

SWARE Buffer

flush 3 pages

tail leaf node

non-overlapping pages

B⁺-tree

non overlapping pages may move

flush non-overlapping pages to tree

bulk load page-by-page if in order

Leaf pages

...

BOSTON UNIVERSITY

DiSC lab

# SWARE Ingestions

Buffer

Zonemap
(min-max)

700-749  750-799  800-849  850-910  870-950  920-991  890-1080  1090-1500

SWARE Buffer

tail leaf node

non-overlapping pages

B$^+$-tree

📌 non overlapping pages may move

📌 flush non-overlapping pages to tree

📌 bulk load page-by-page if in order

Leaf pages

...

# SWARE Ingestions

move & sort remaining entries

Buffer

SWARE Buffer

Zonemap
(min-max)

850-910    915-970    974-1030    1032-1088    1090-1500

tail leaf node

non-overlapping pages

B⁺-tree

non overlapping pages may move

flush non-overlapping pages to tree

bulk load page-by-page if in order

Leaf pages                 ...

BOSTON UNIVERSITY

DiSC lab

# SWARE Ingestions

Buffer

Zonemap
(min-max)



update non-overlapping pages

SWARE Buffer

850-910    915-970    974-1030    1032-1088    1090-1500

tail leaf node

non-overlapping pages

B⁺-tree

non overlapping pages may move

flush non-overlapping pages to tree

bulk load page-by-page if in order

Leaf pages          . . .

BOSTON
UNIVERSITY

DiSC lab

# SWARE Ingestions

Buffer

Zonemap
(min-max)

SWARE Buffer

850-910    915-970    974-1030    1032-1088    1090-1500

tail leaf node

non-overlapping pages

fully sorted pages

B$^+$-tree

📌 non overlapping pages may move

📌 flush non-overlapping pages to tree

📌 bulk load page-by-page if in order

Leaf pages          ...

# How do lookups work?

# SWARE Lookups

Global Bloom filter

Zonemap (min-max)

850-899  900-950  951-1000  1111-1168  1170-1500  1002-1200  1010-1800  1070-1999

SWARE Buffer

**Per-page Bloom filters**

**Buffer**

tail leaf node

non-overlapping pages

fully sorted pages

sorted section uses faster interpolation search

Global BF. helps skip buffer probe

Per-page BFs + Zonemaps eliminate page-scans

B$^+$-tree

Leaf pages   ...

BOSTON UNIVERSITY

DiSC lab

# Experimental Evaluation

**System Setup:**

- Intel Xeon Gold 5230

- 2.1GHZ processor w. 20 cores

- 384GB RAM, 28MB L3 cache

**Index Setup:**

- Buffer = 40MB; flush <= 50%

- BFs = 10 BPK; Murmur Hash

- Split at 80%

$B^+$-tree design inspired by STX::B-tree can also work as $B^\varepsilon$-tree

# Evaluating SWARE Under Varying Sortedness
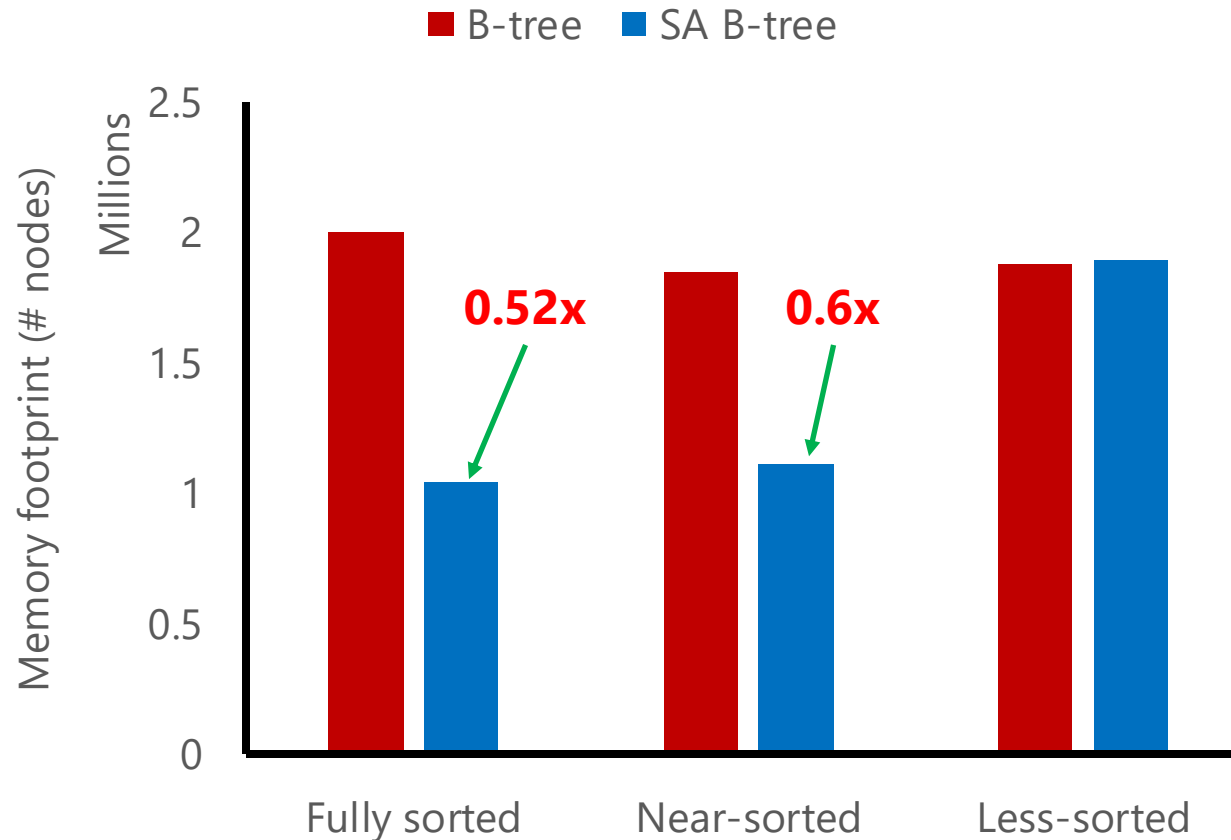
# Raw Ingestion Performance



ingestion latency reduced between 27-90%
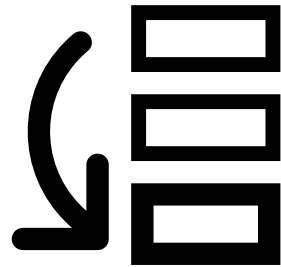
bulk loading maximized with
high data sortedness

# SWARE Improves Space Utilization



**increased fill/split factor helps reduce memory footprint**

# Summarizing SWARE [ICDE 2023]



intelligent
buffering

+

opportunistic
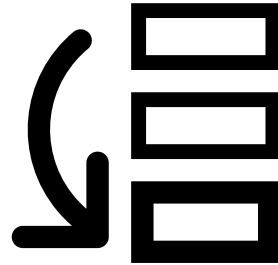bulk loading

Improves performance by
exploiting sortedness

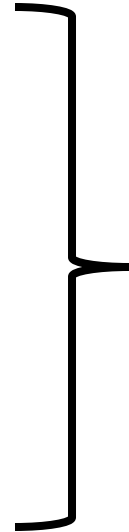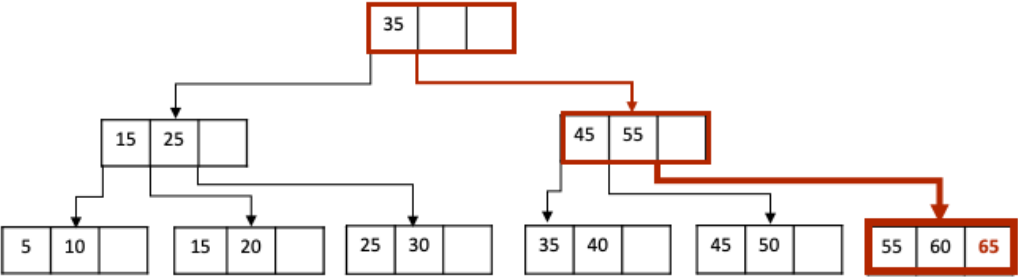**Any downsides to wider applicability?**

# Summarizing SWARE [ICDE 2023]



intelligent buffering + opportunistic bulk loading } Improves performance by exploiting sortedness

**Increases Complexity in Design!**

# Can we achieve fast ingestions **without** buffering?

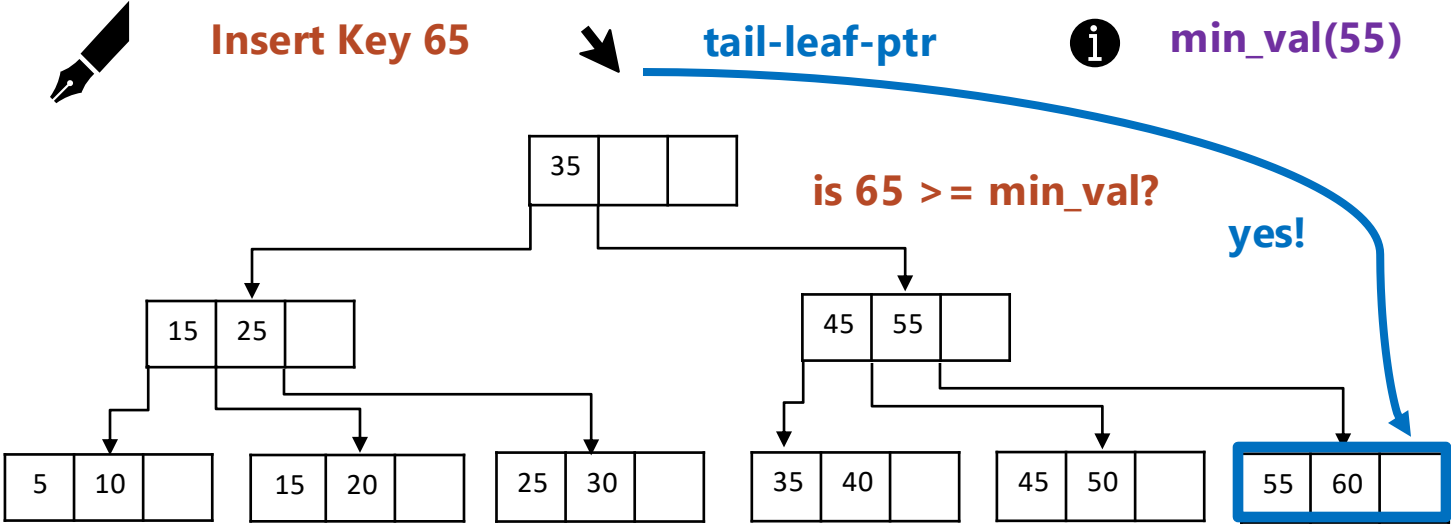# Inserting to the Tail-leaf (PostgreSQL & MySQL)
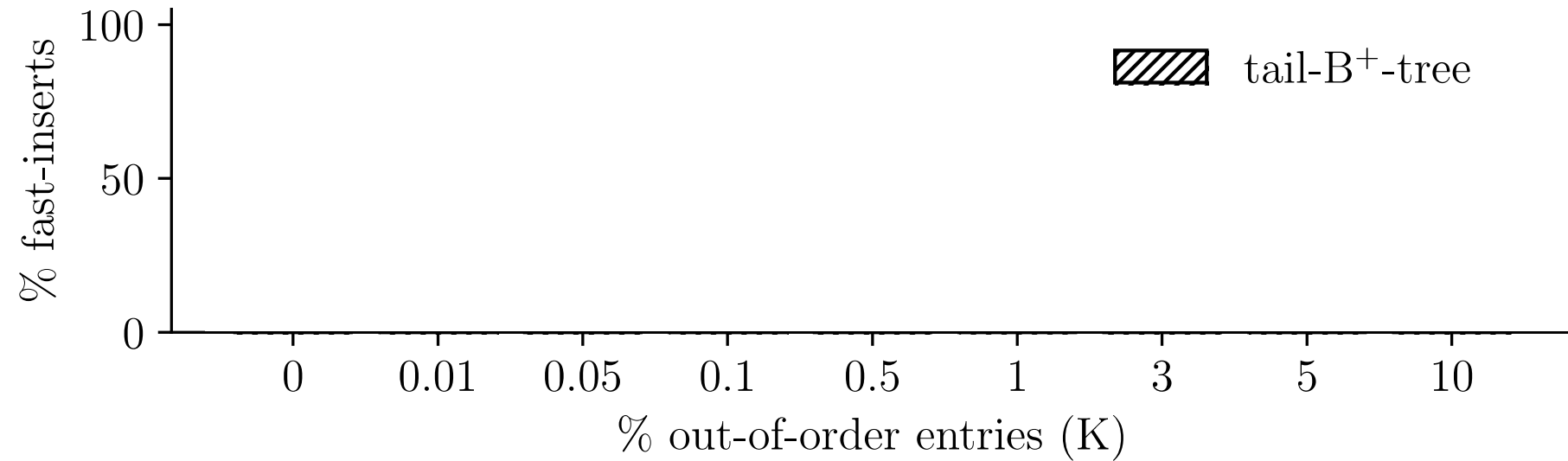


Normal Insertion (top-insert)

Tail-leaf Insertion

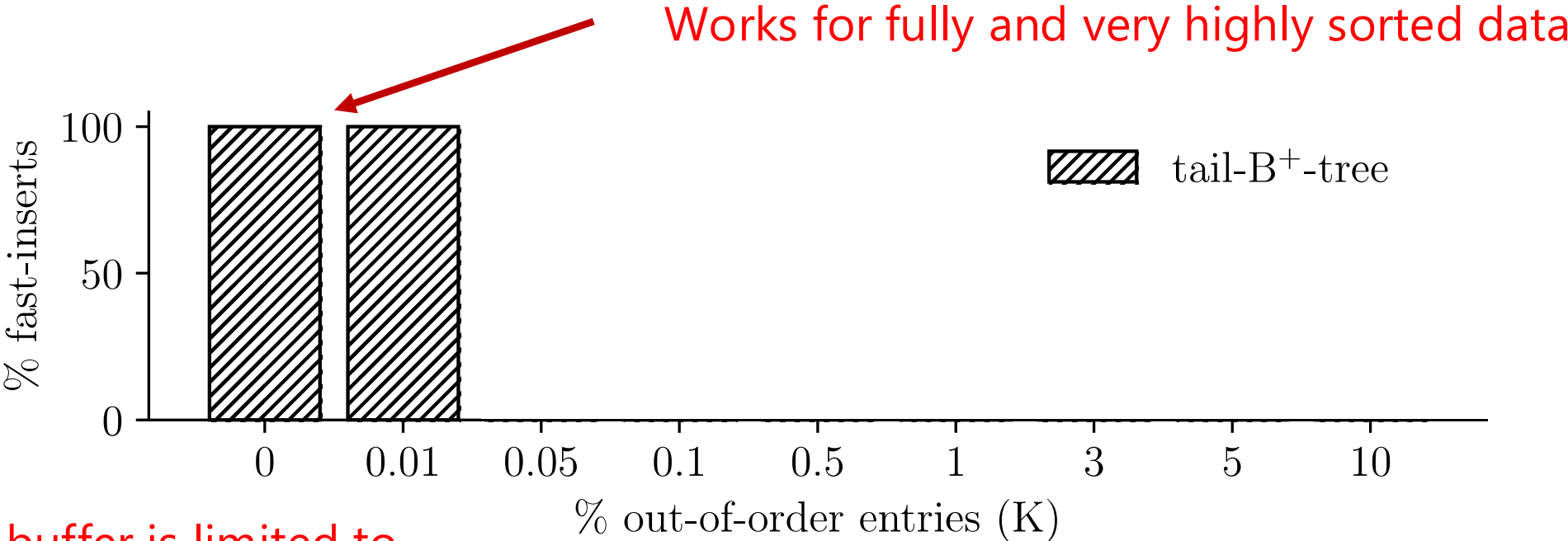# Is the tail-leaf optimization the solution?

# Does Tail-leaf Insertion Work?



(plot) y-axis: % fast-inserts (0, 50, 100); x-axis: % out-of-order entries (K) with values 0, 0.01, 0.05, 0.1, 0.5, 1, 3, 5, 10. Legend: tail-$B^+$-tree

BOSTON UNIVERSITY

DiSC lab

# Does Tail-leaf Insertion Always Work?



Works for fully and very highly sorted data

Tail-leaf's buffer is limited to leaf node!

# Does Tail-leaf Insertion Always Work?



Works for fully and very highly sorted data

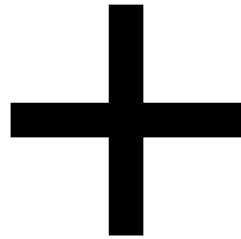Tail-leaf's buffer is limited to leaf node!

Degrades very quickly

**However, tail-leaf points us to the right direction…**
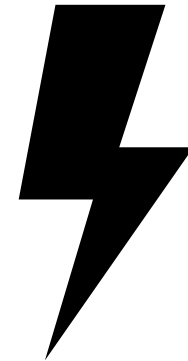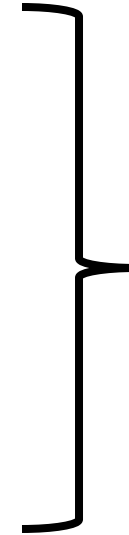
# Key Idea – Predicting the Ordered LEaf (POLE)
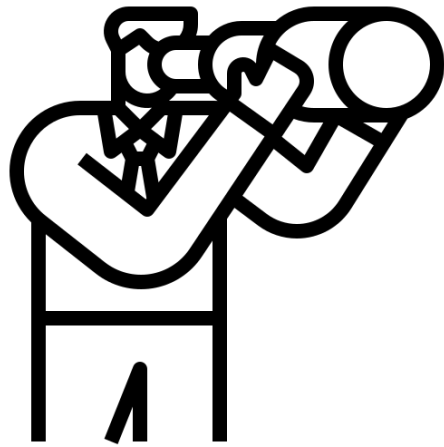


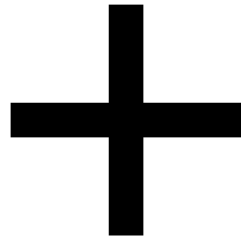Sortedness-aware predictor

Leaf appends

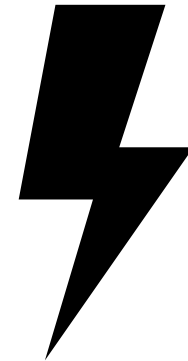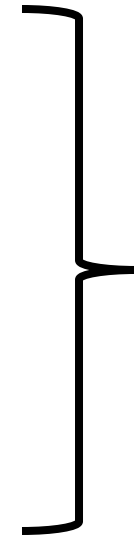Fast ingestion

# Key Idea – Predicting the Ordered LEaf (POLE)

it could be *any* node!



Sortedness-aware predictor

+

Leaf appends

Fast ingestion

BOSTON UNIVERSITY

DiSC lab

# Insertions in Steady-State

Insert *(x, v)*

top-insert

if **x** is in **pole** range & fits



B⁺-tree

*pole*   ...   *tail*

# When Pole Splits

**Predict using _IKR (In-order Key estimatoR)_**

$$x = q + \left( \frac{q - p}{pole\_prev_{size}} \right) \cdot pole_{size} \cdot (1.5)$$

density between two non-outliers

B⁺-tree

| p | | q | | r | |

_pole_prev_   _current pole_   _newly created node_

BOSTON UNIVERSITY

DiSC lab

# When Pole Splits

if r > x, new node has outliers



B⁺-tree

*p* | | *q* | | *r* |

*pole_prev*      *current pole*      *newly created node*

**pole** stays as is

**Legend**
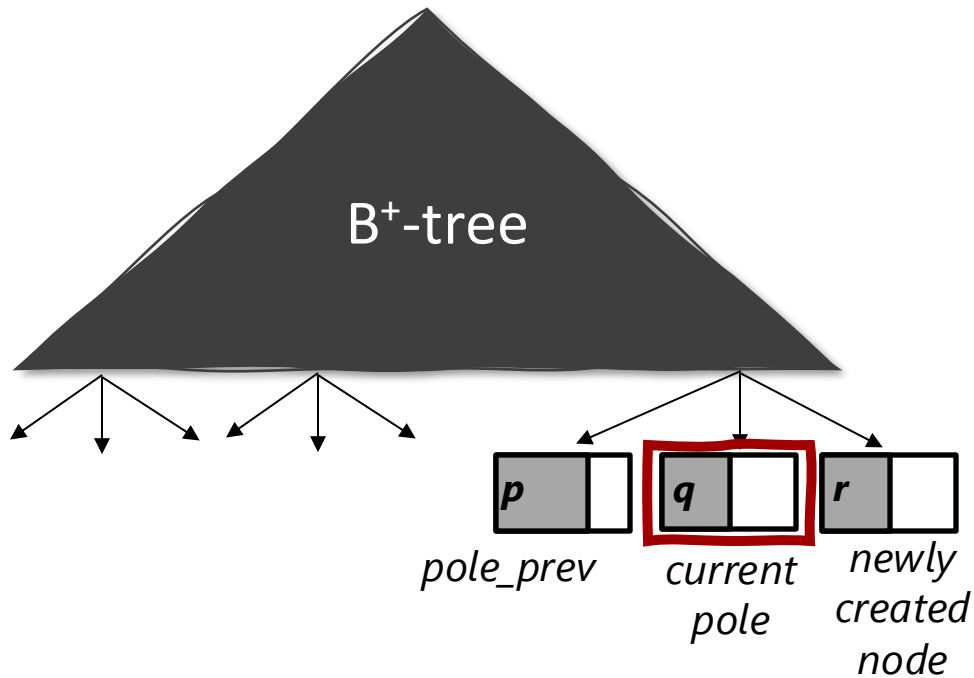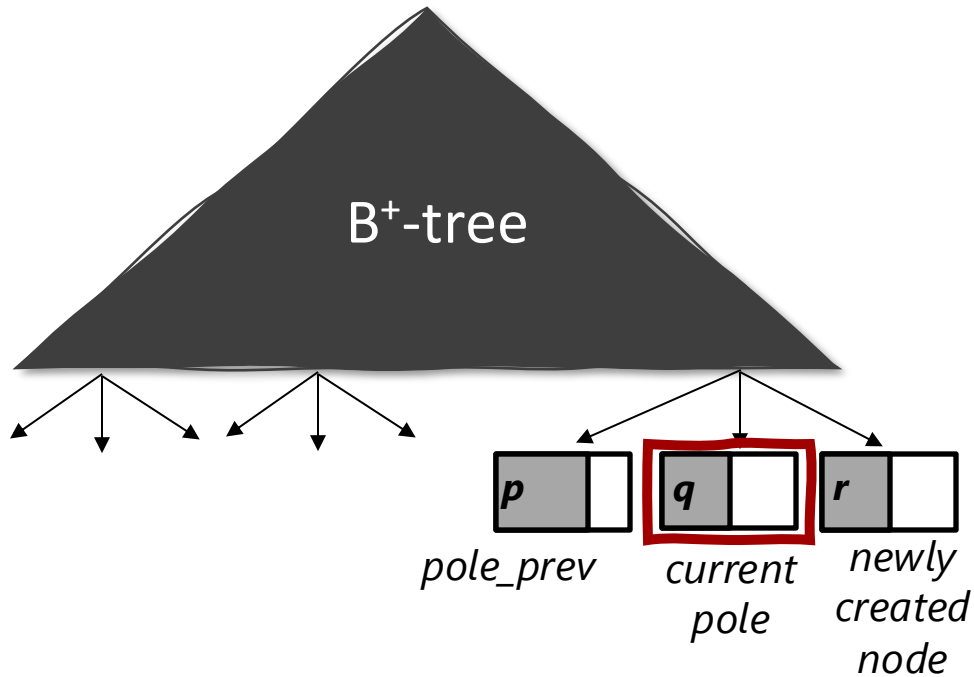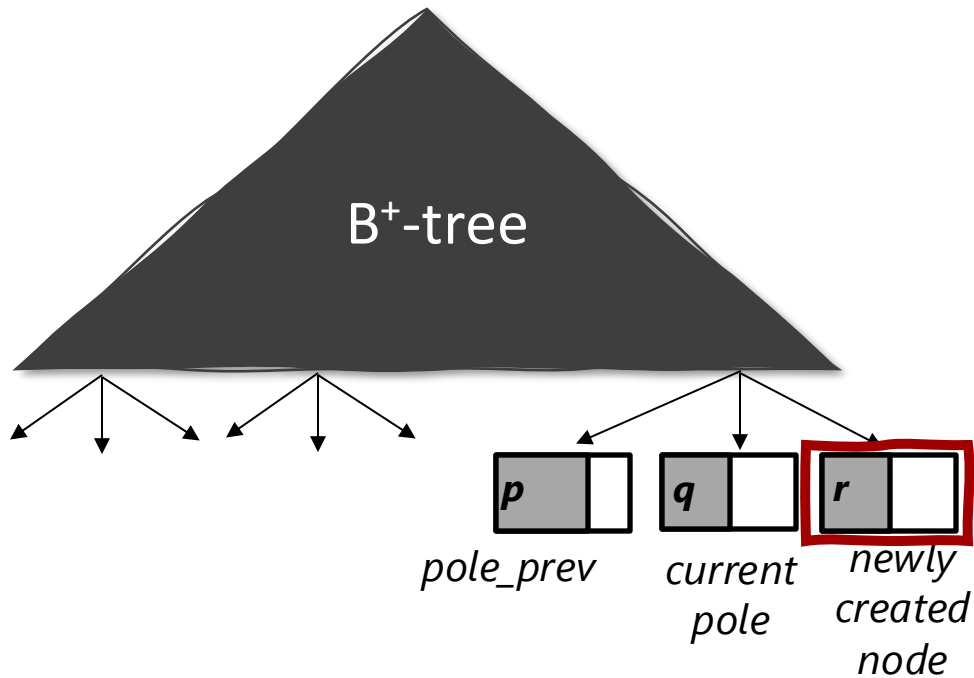**p** = smallest entry in node previous to *pole;*
**q** = smallest entry in *pole*
**r** = smallest entry in newly created node
☐ = pointer to *pole* node
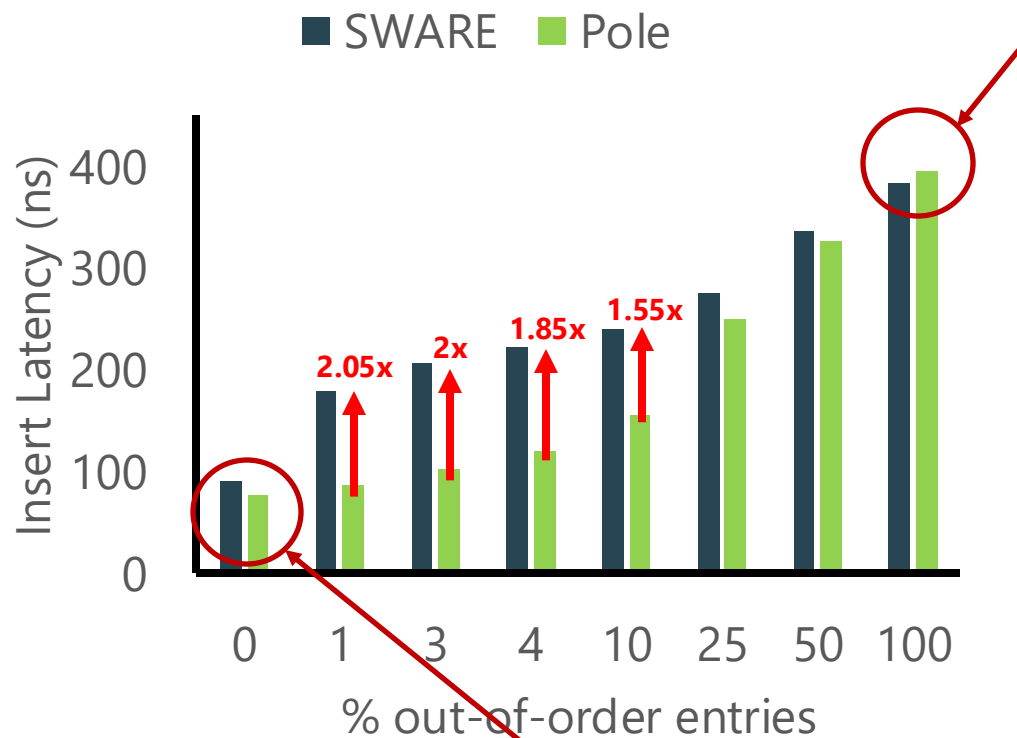
**Predict using *IKR (In-order key estimatoR)***

$$x = q + \left( \frac{q - p}{pole\_prev_{size}} \right) \cdot pole_{size} \cdot (1.5)$$

density between two non-outliers

# When Pole Splits

if r <= x, new node has at least one non-outlier value

**Legend**
**$p$** = smallest entry in node previous to *pole*;
**$q$** = smallest entry in *pole*
**$r$** = smallest entry in newly created node
☐ = pointer to *pole* node

B⁺-tree

p       q       r

*pole_prev*   *current pole*   *newly created node*

**Predict using *IKR (In-order key estimatoR)***

$$x = q + \left(\frac{q - p}{pole\_prev_{size}}\right) \cdot pole_{size} \cdot (1.5)$$

density between two non-outliers

Update ***pole*** to newly created node from split

BOSTON UNIVERSITY

DiSC lab

# Comparing with SWARE



SWARE   Pole

Buffer pays off: some vs. none fast ingestion

**up to 2.05x faster**

minimal metadata ✓

avoids SWARE buffer management ✓

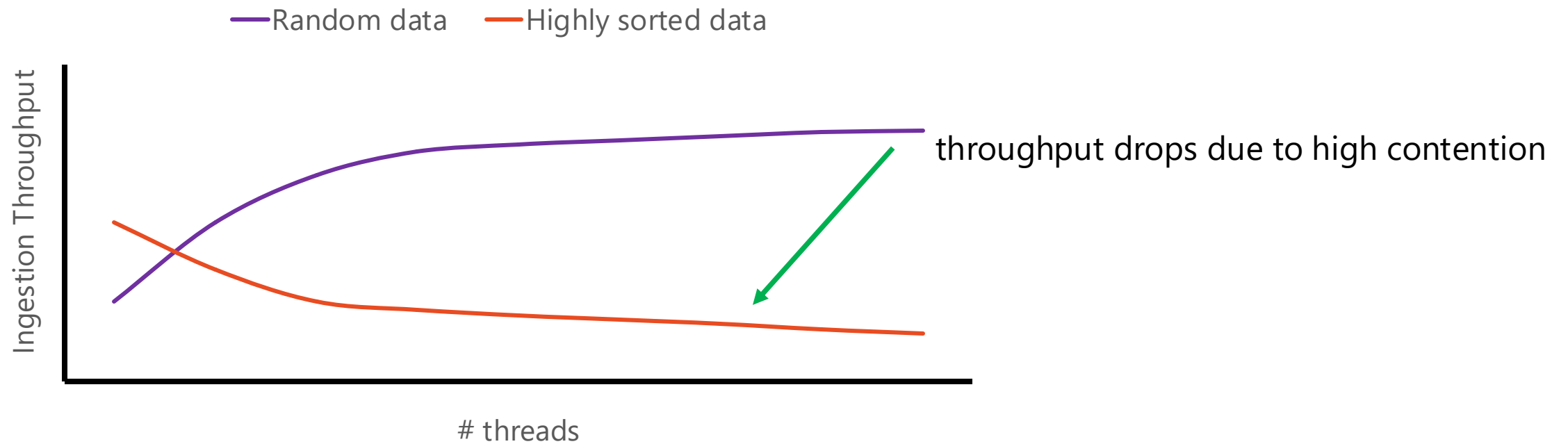buffer helps: full bulk loading
Pole is still faster!
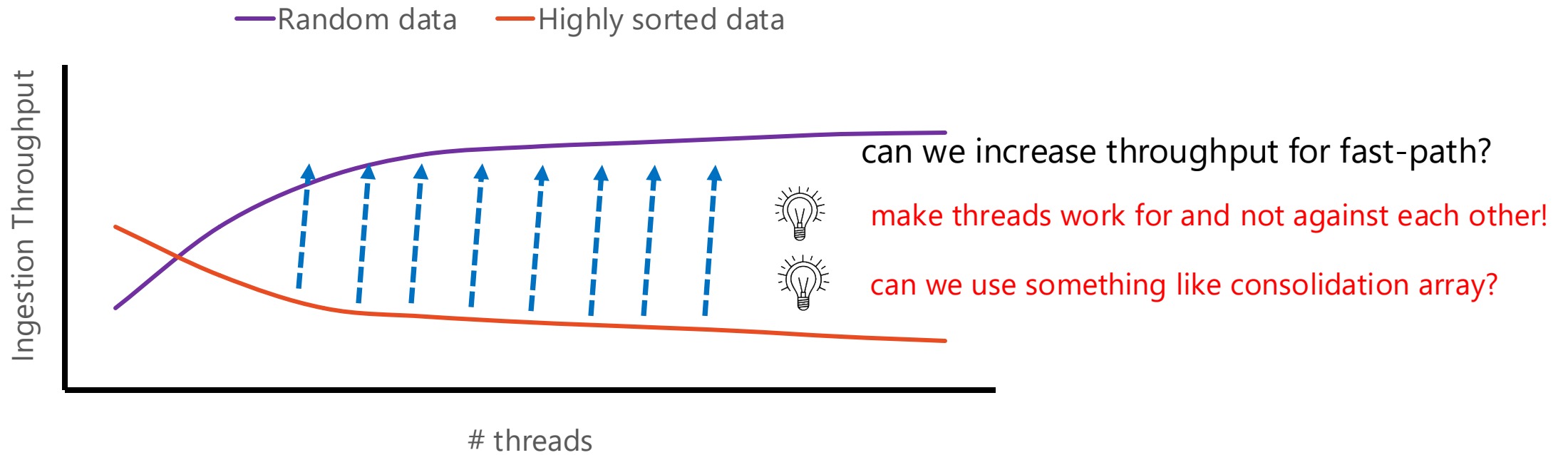
# Comparing with SWARE



full bulk loading ⇒ smaller tree

up to 29% faster for point lookups

No buffering ⇒ no read overhead!

# Future Work - Concurrency in Fast Path
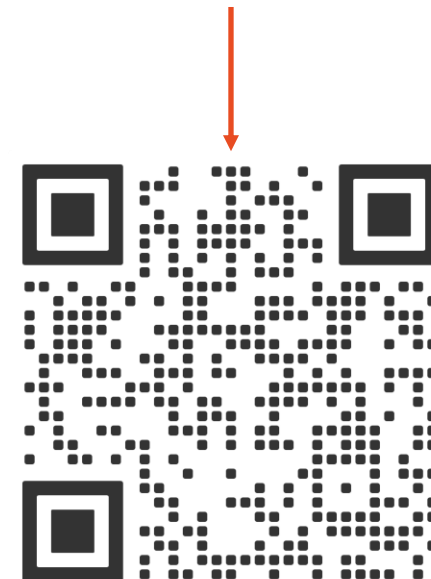
# Future Work - Concurrency in Fast Path

# Summary

Identify "sortedness" as a resource

Classical indexes do not exploit sortedness by design!

SWARE paradigm & Pole optimization optimize for sortedness

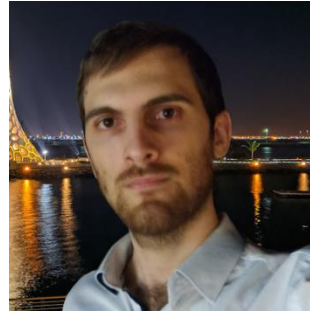Further research required for learned indexes + joins

BOSTON UNIVERSITY
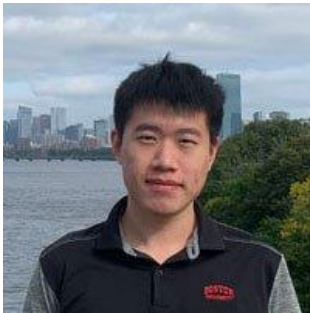
DiSC lab

# The Team


Aneesh Raman


Konstantinos
Karatsenidis


Andy Huynh


Jinqi Lu


Shaolin Xie


Subhadeep  Sarkar


Matthaios Olma