

CS 561: Data Systems Architectures

Introduction to Indexing:

Trees, Tries, Hashing, Bitmaps: The Whole Design Space

Dr. Subhadeep Sarkar

<https://bu-disc.github.io/CS561/>

A few reminders

- A) Choose your **class project**. (40% of your grades)
- B) Project proposals **due on 03/12** in groups of 2-3.
- C) Proposal format:
 - (i) Project background
 - (ii) Problem statement
 - (iii) Project objectives
 - (iv) Project timeline



What is an index?

Auxiliary structure to quickly find rows based on arbitrary attribute

Special form of <key, value>



indexed attribute



position/location/rowID/primary key/...

What are the possible *index designs*?

	Data Organization	Comments
B+ Trees	Sorted & partitioned	Partition <i>k-ways</i> recursively
LSM Trees	Insertion & Sorted	Optimizes <i>insertion</i>
Radix Trees	Radix	Partition using the <i>key radix</i> representation
Hash Indexes	Hash	Partition by <i>hashing the key</i>
Bitmap Indexes	None	Succinctly represent <i>all rows with a key</i>
Scan Accelerators	None	Metadata to <i>skip accesses</i>

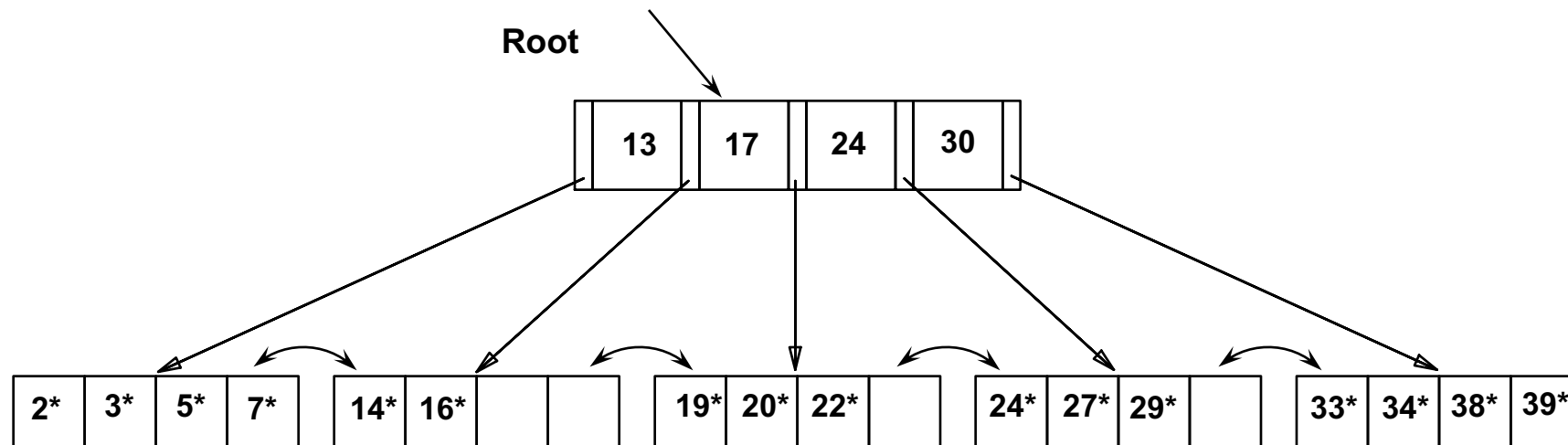
What are the possible *index designs*?

	Point Queries	Short Range Queries	Long Range Queries	Data Skew	Updates
B+ Trees					
LSM Trees					
Radix Trees					
Hash Indexes					
Bitmap Indexes					
Scan Accelerators					

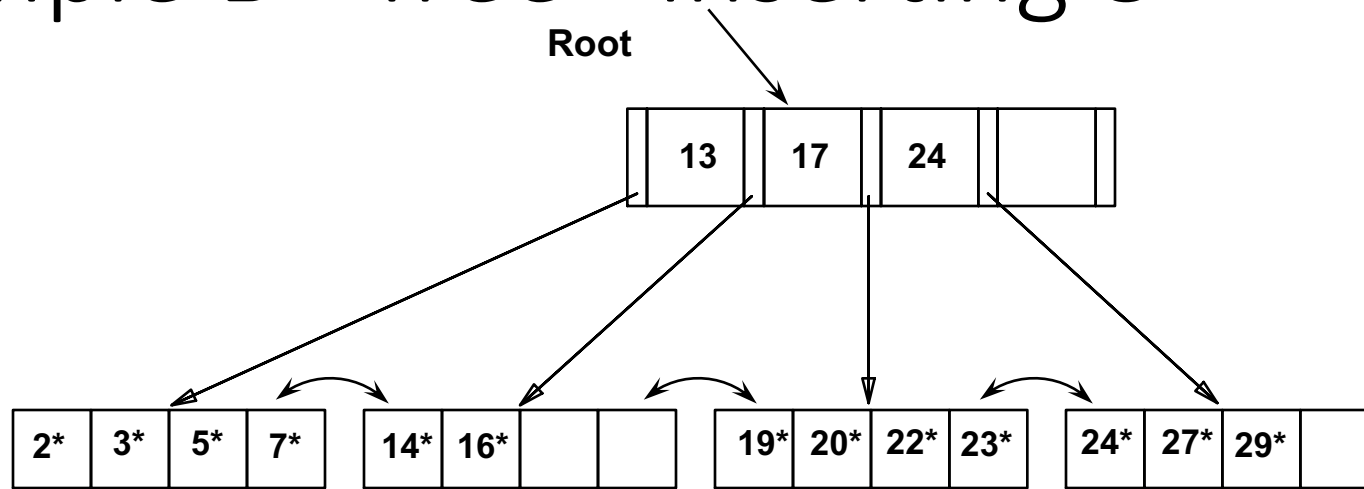
B+ Trees

Search begins at root, and key comparisons direct it to a leaf.

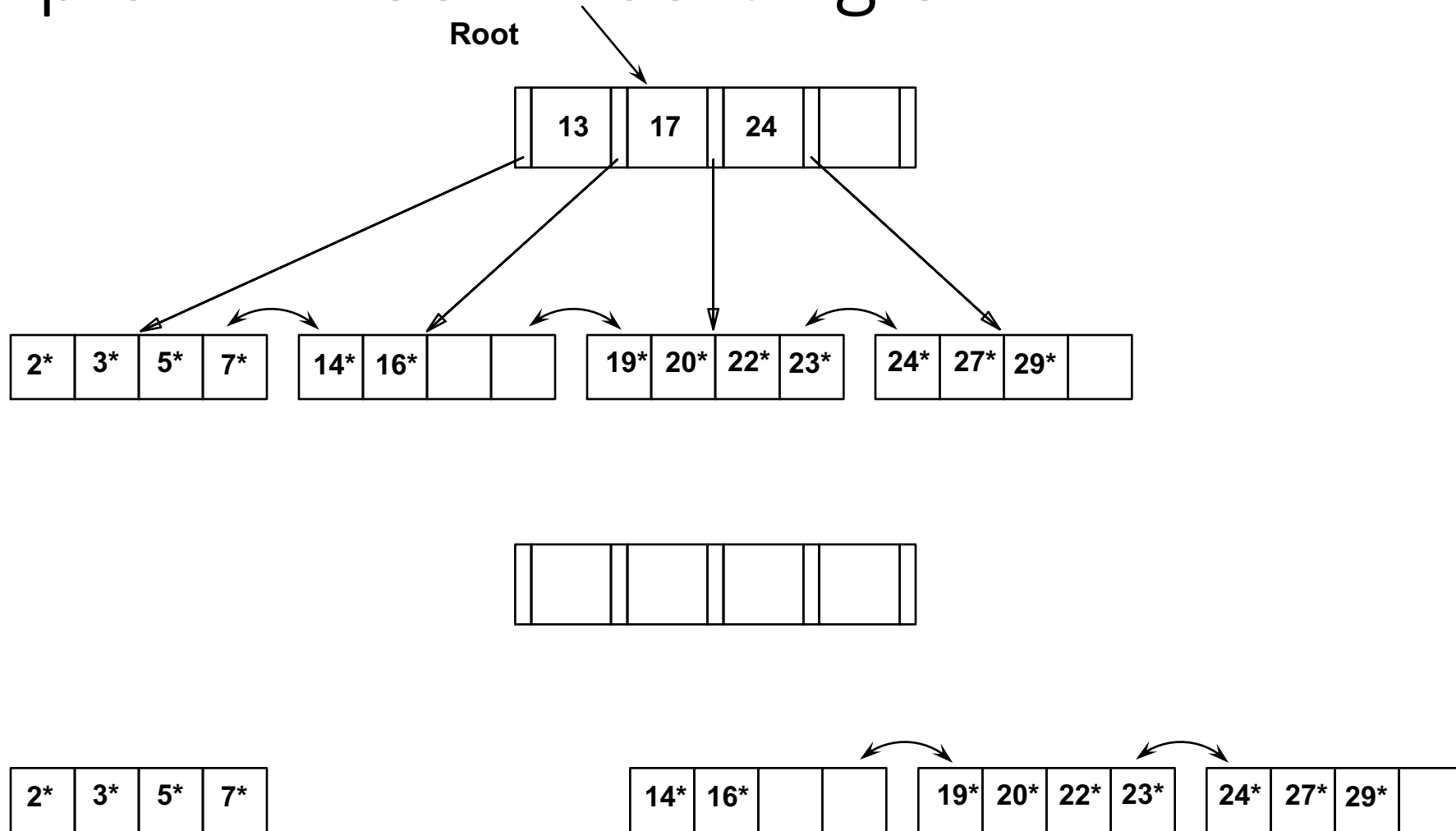
Search for 5^* , 15^* , all data entries $\geq 24^*$...



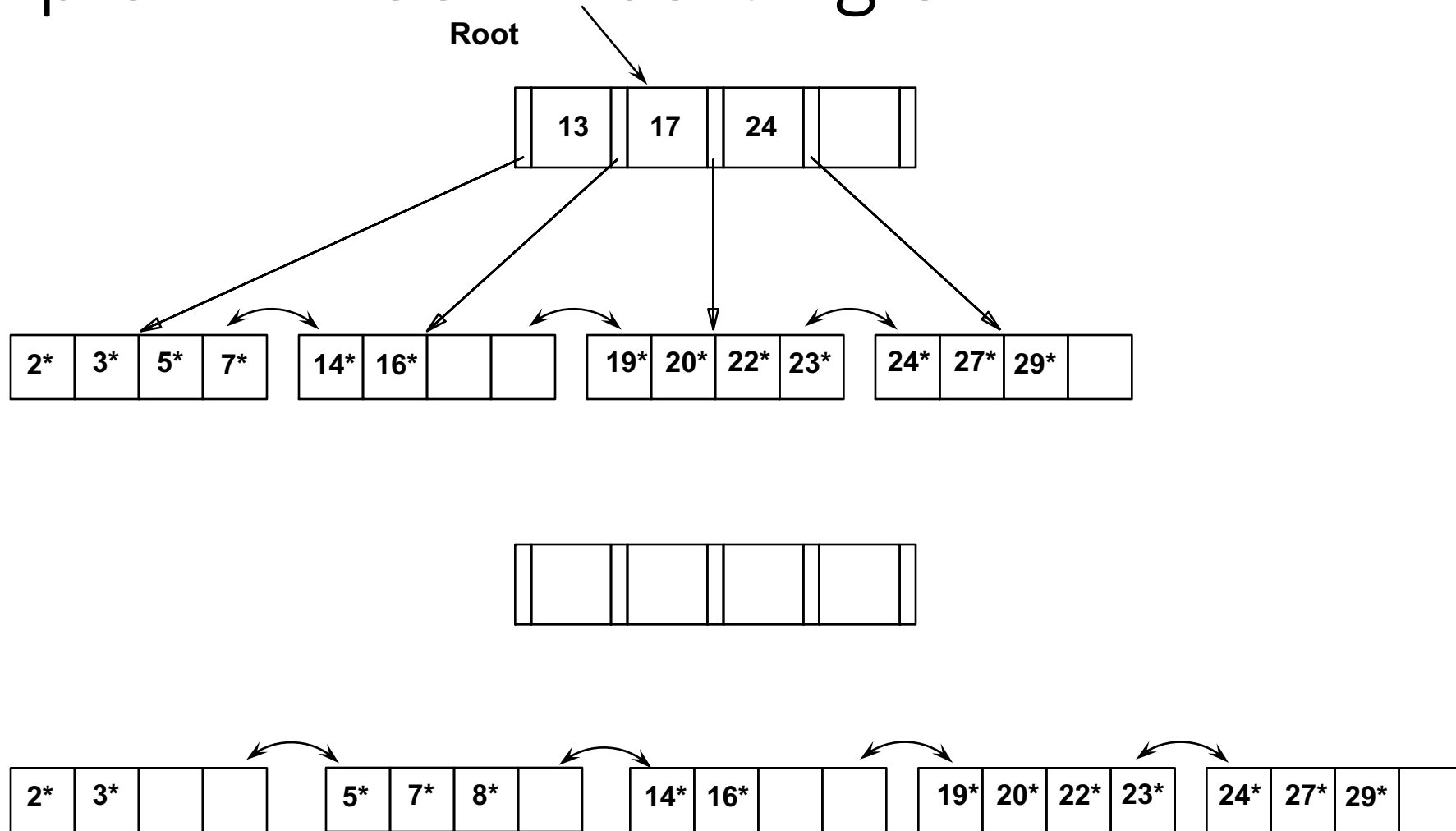
Example B+ Tree - Inserting 8*



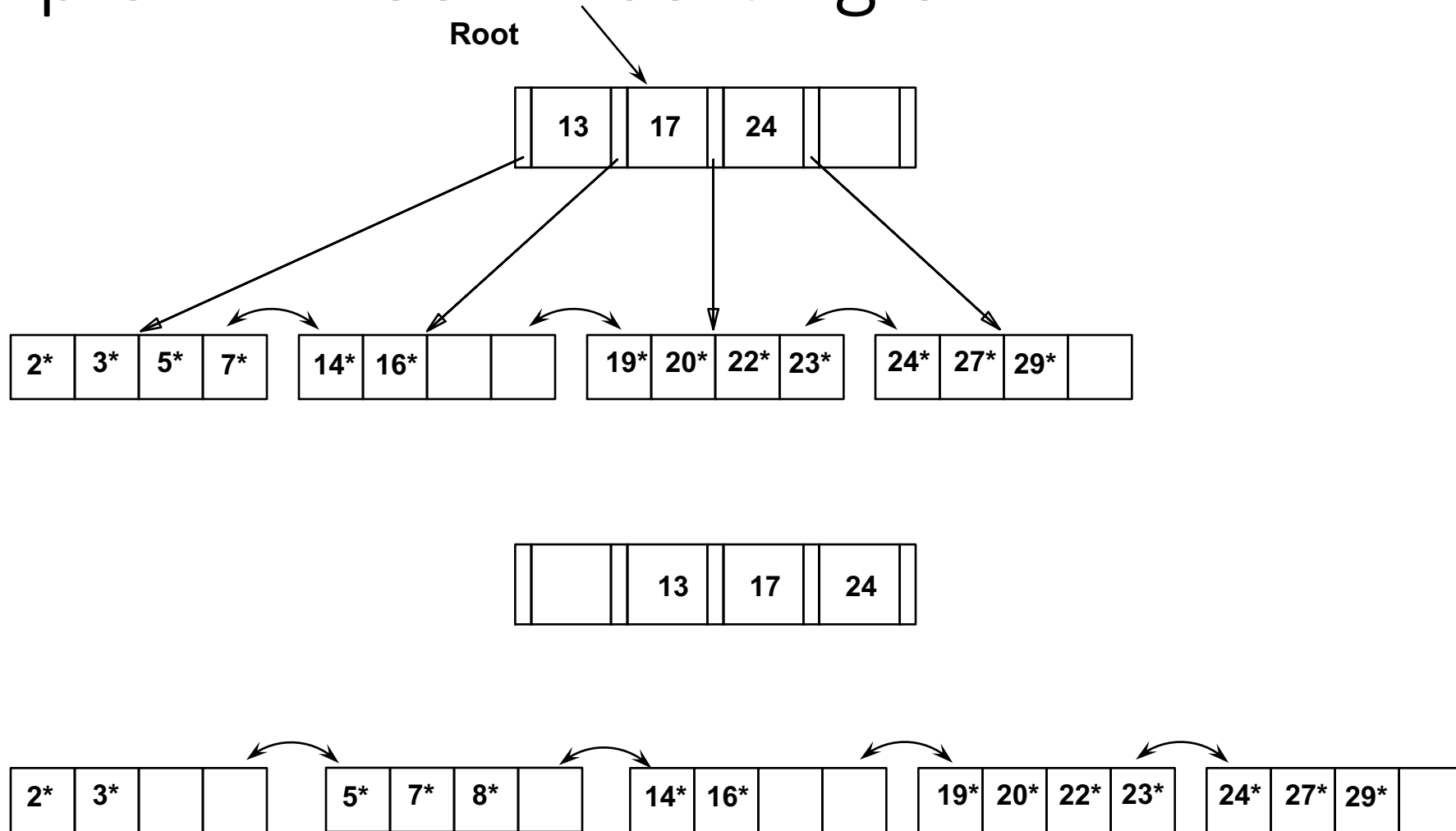
Example B+ Tree - Inserting 8*



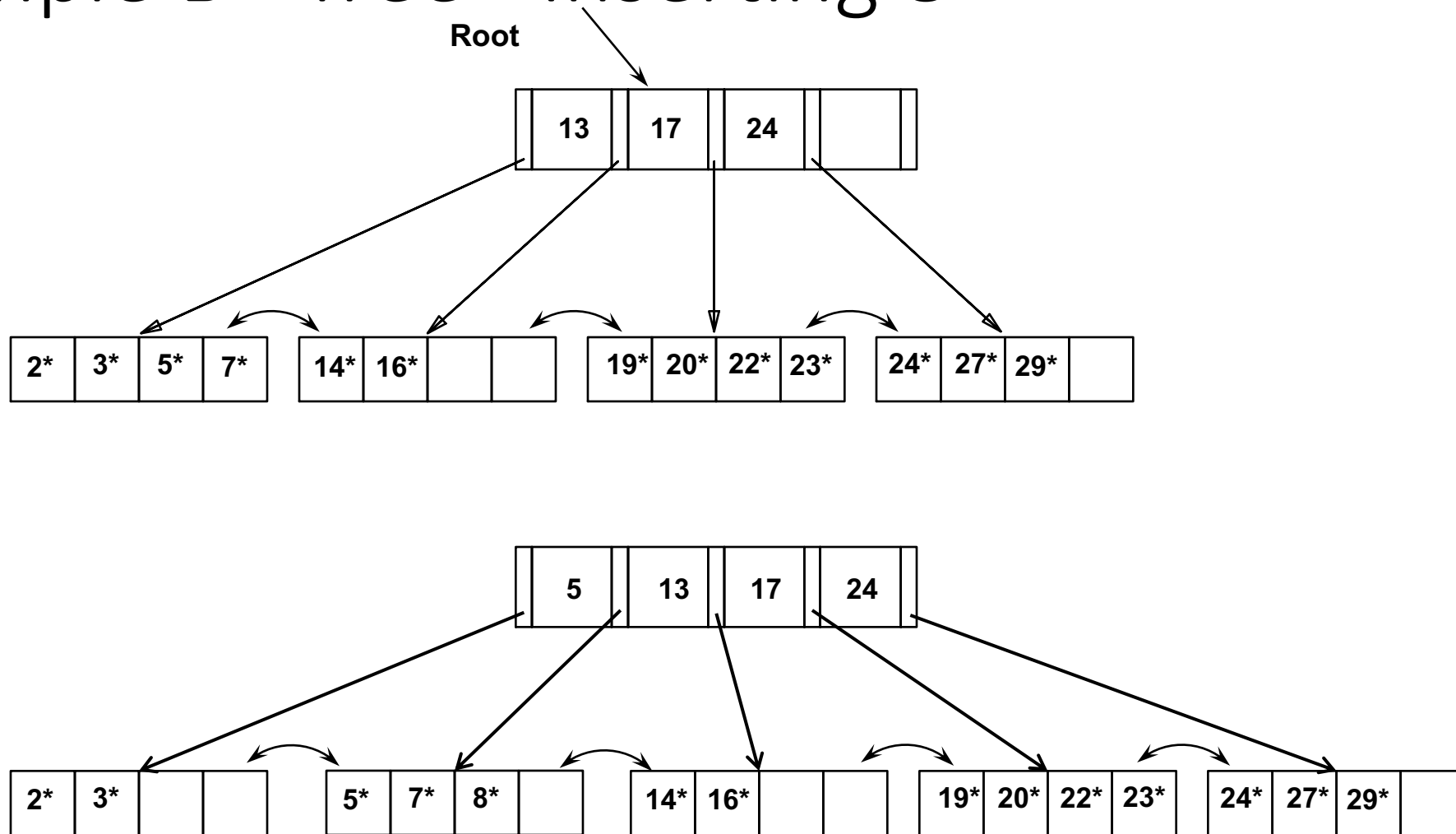
Example B+ Tree - Inserting 8*



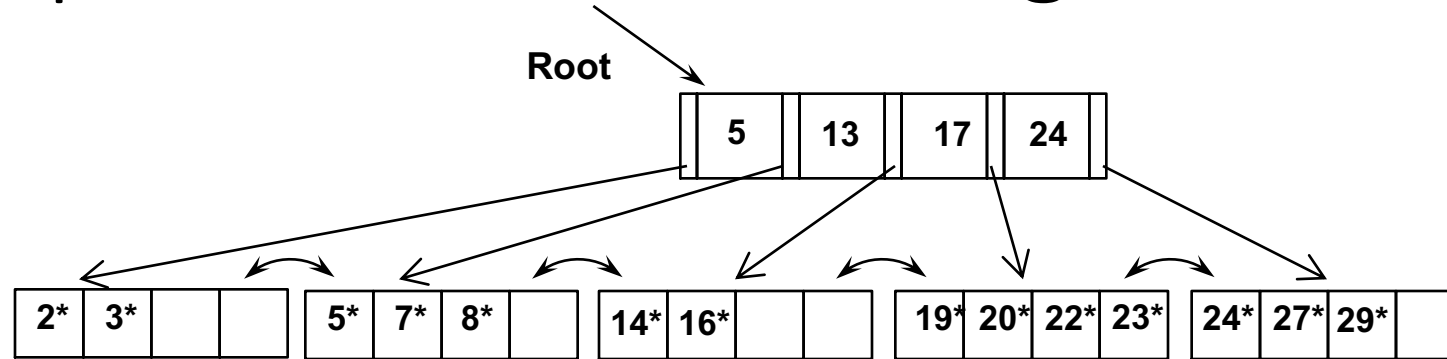
Example B+ Tree - Inserting 8*



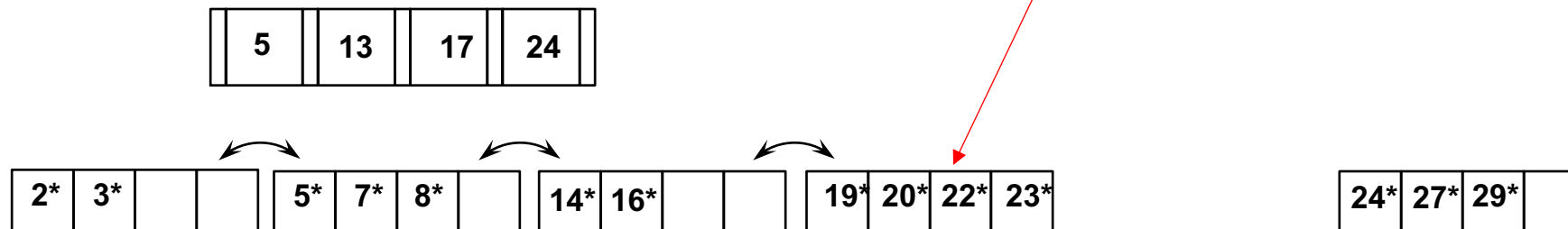
Example B+ Tree - Inserting 8*



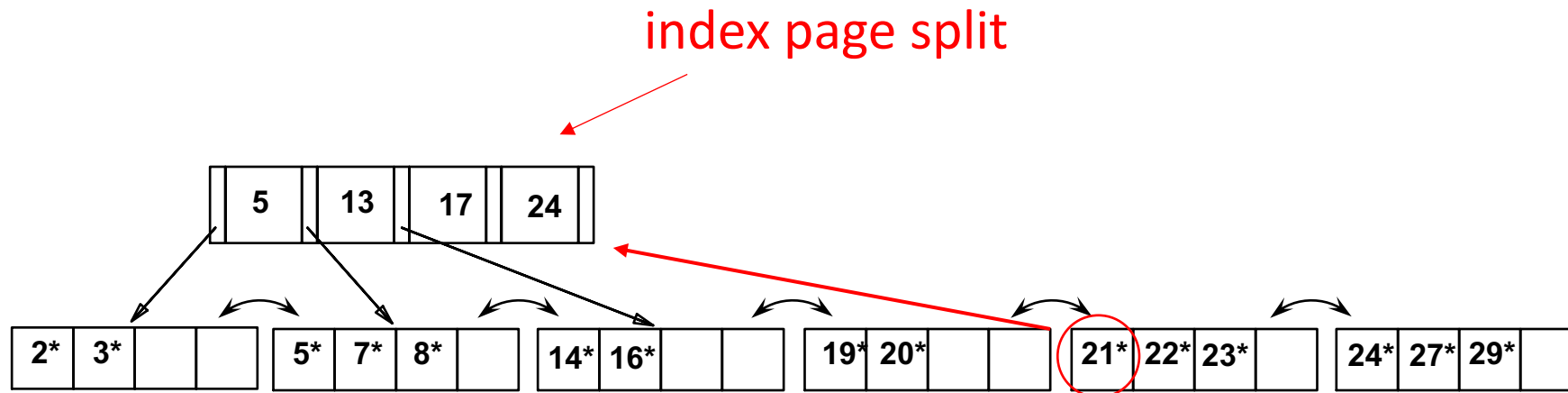
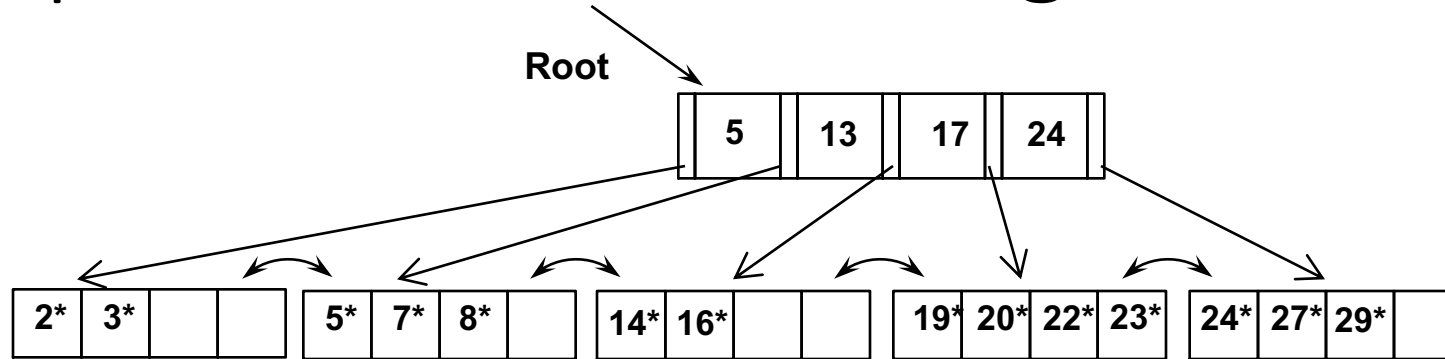
Example B+ Tree - Inserting 21*



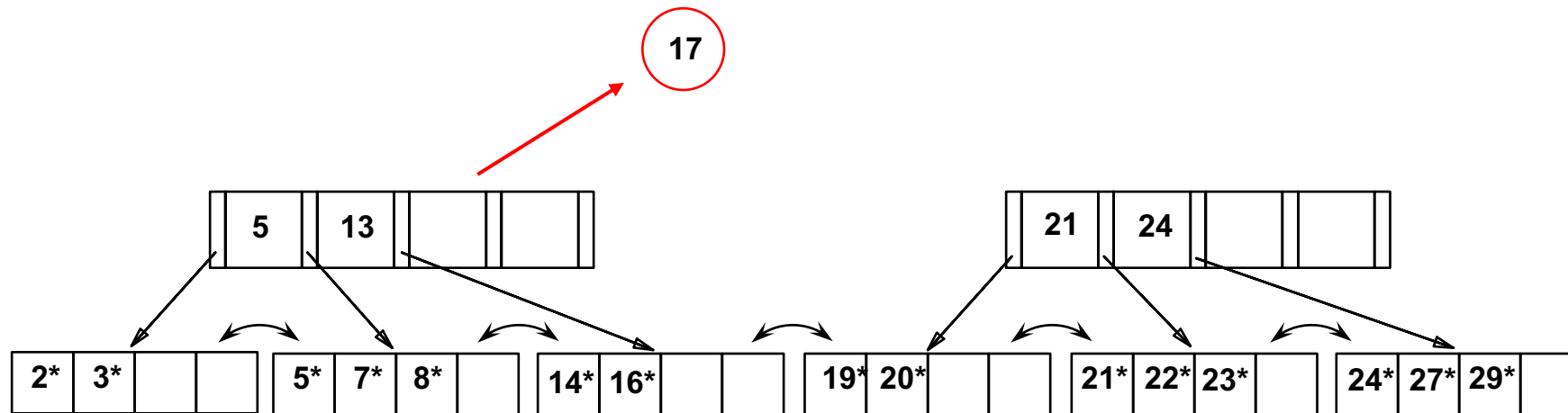
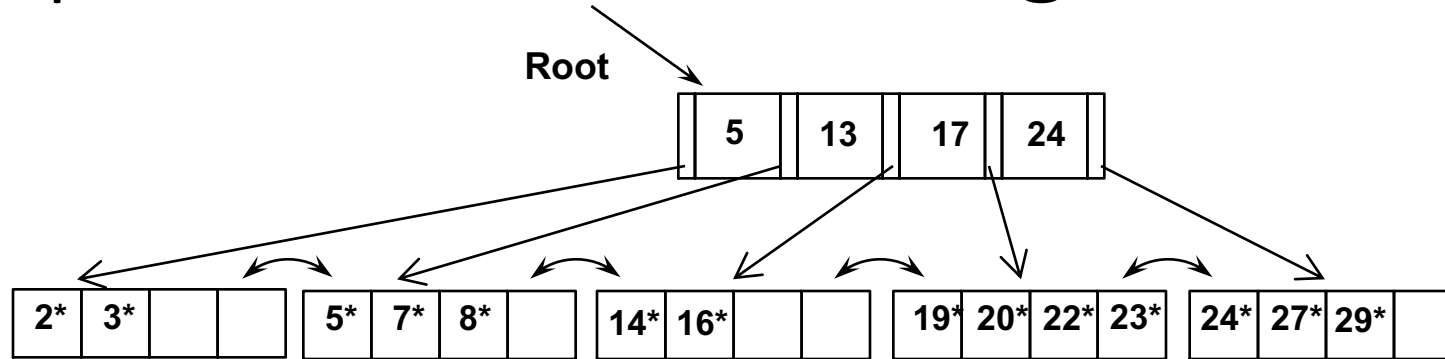
data page split



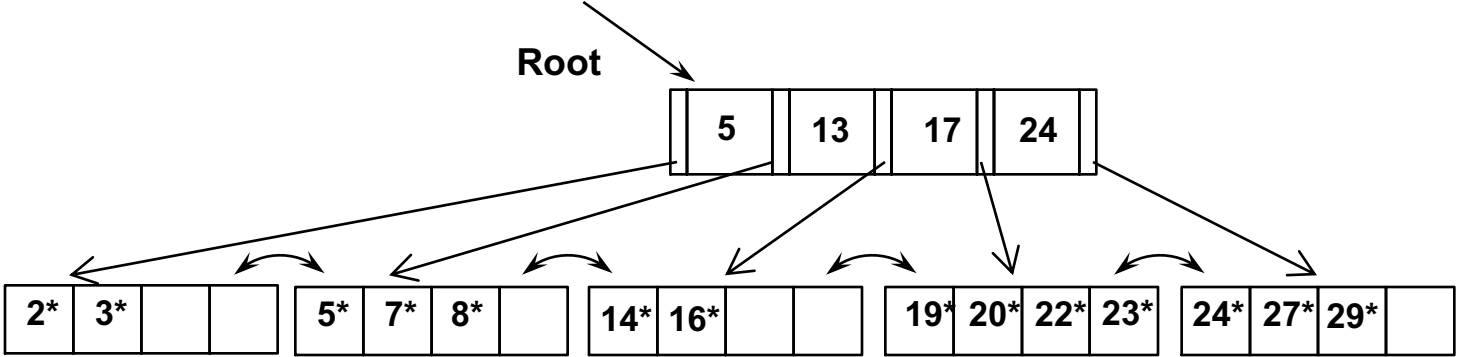
Example B+ Tree - Inserting 21*



Example B+ Tree - Inserting 21*



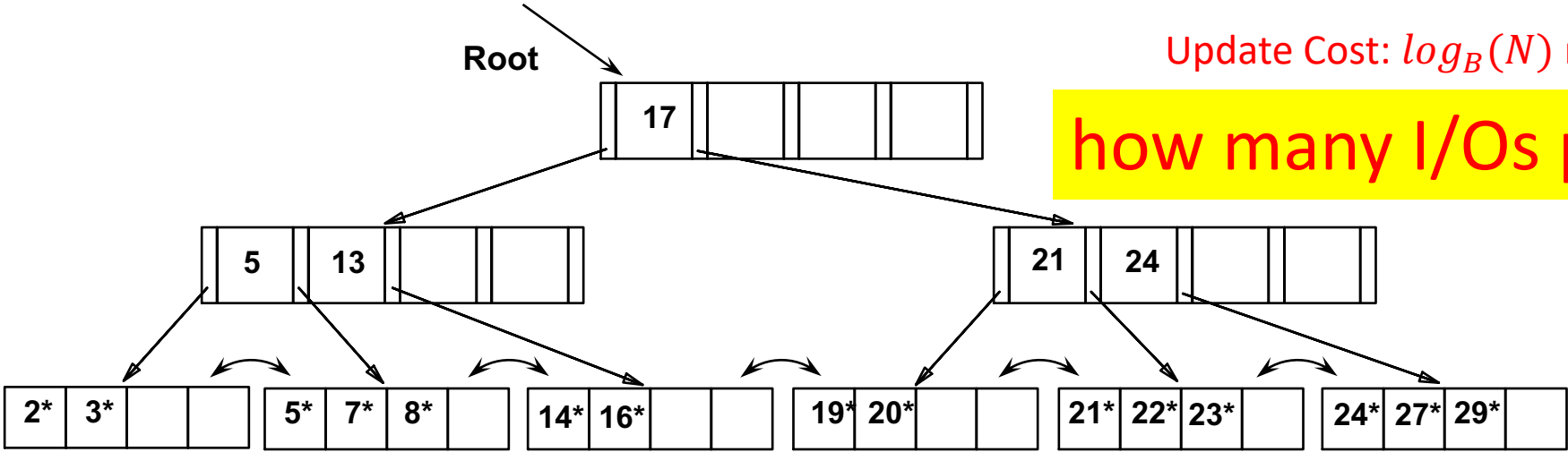
Example B+ Tree - Inserting 21*



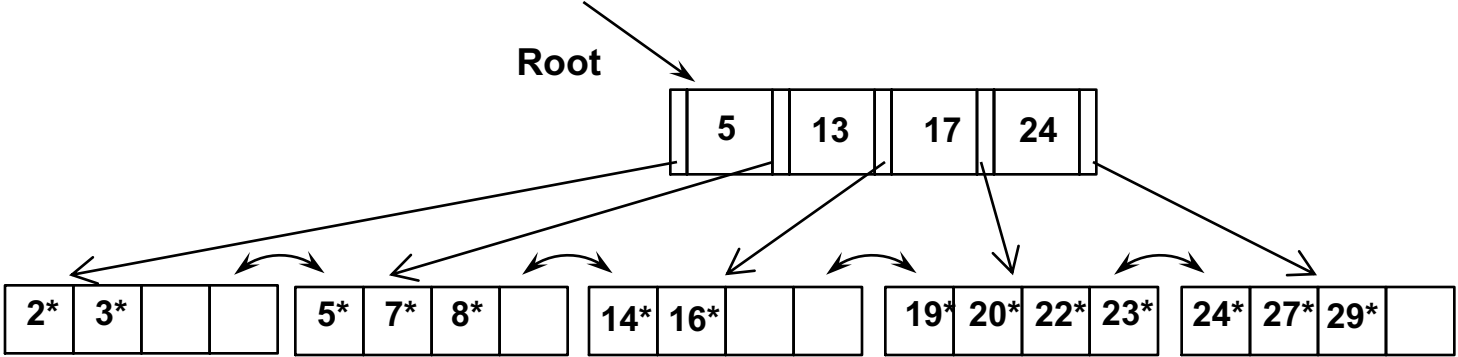
Read Cost: $\log_B(N)$

Update Cost: $\log_B(N)$ reads

how many I/Os per insert?



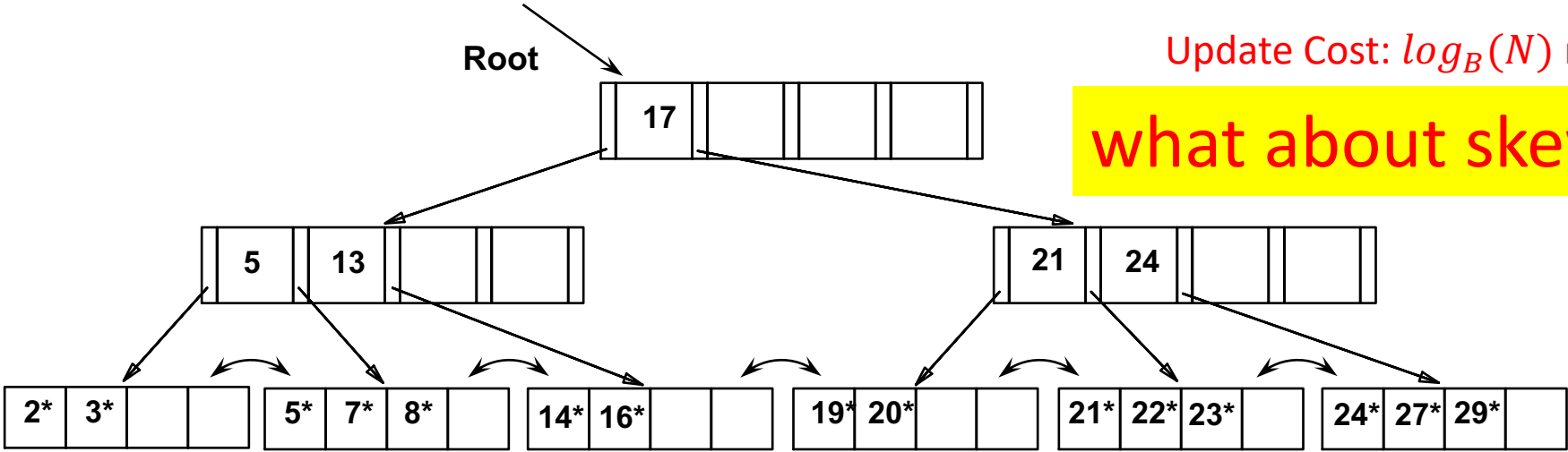
Example B+ Tree - Inserting 21*



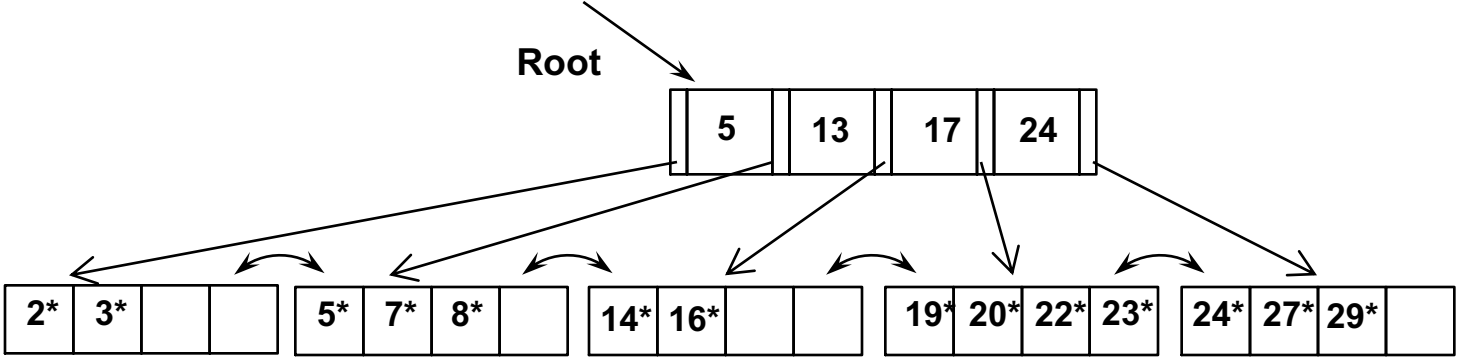
Read Cost: $\log_B(N)$

Update Cost: $\log_B(N)$ reads

what about skewed data?



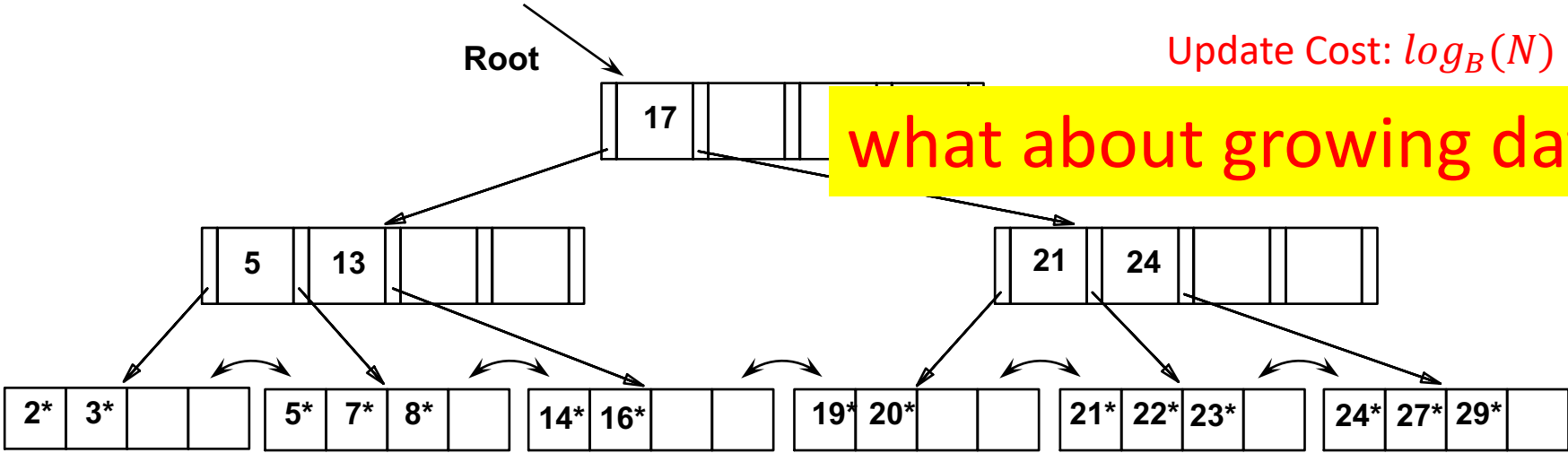
Example B+ Tree - Inserting 21*



Read Cost: $\log_B(N)$

Update Cost: $\log_B(N)$ reads

what about growing dataset size?

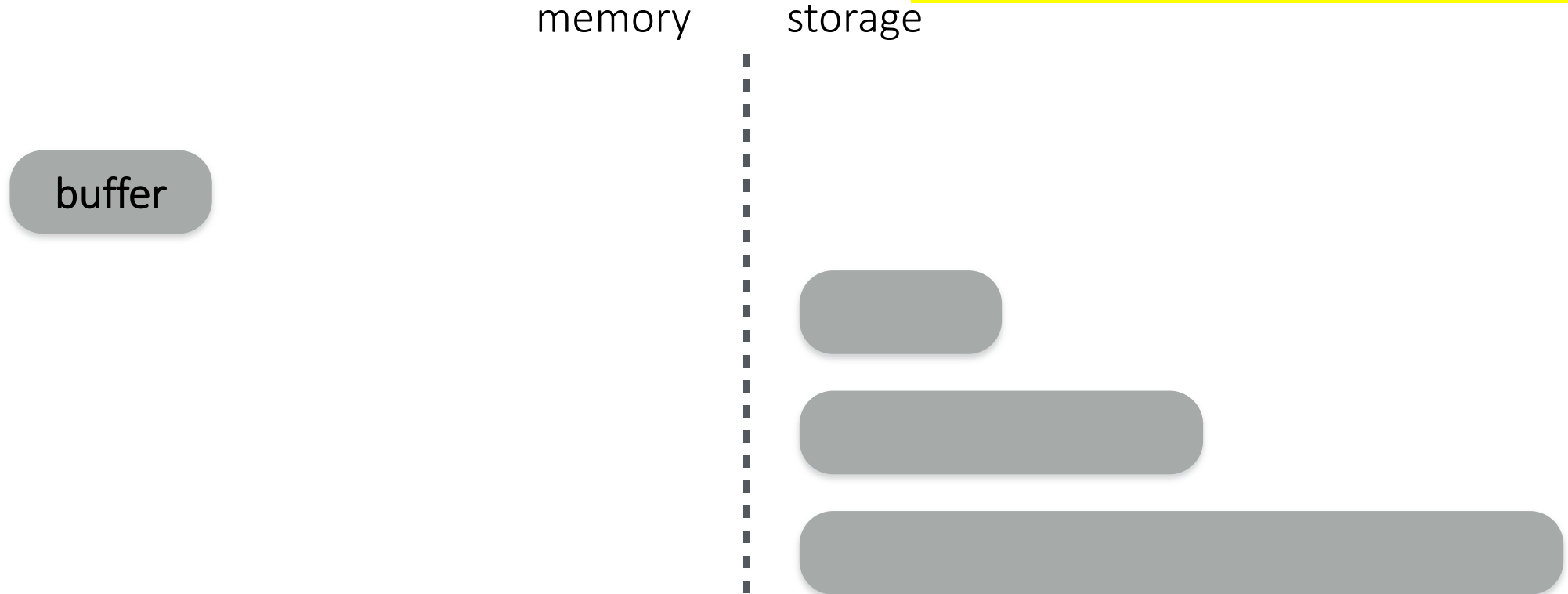


What are the possible index designs?

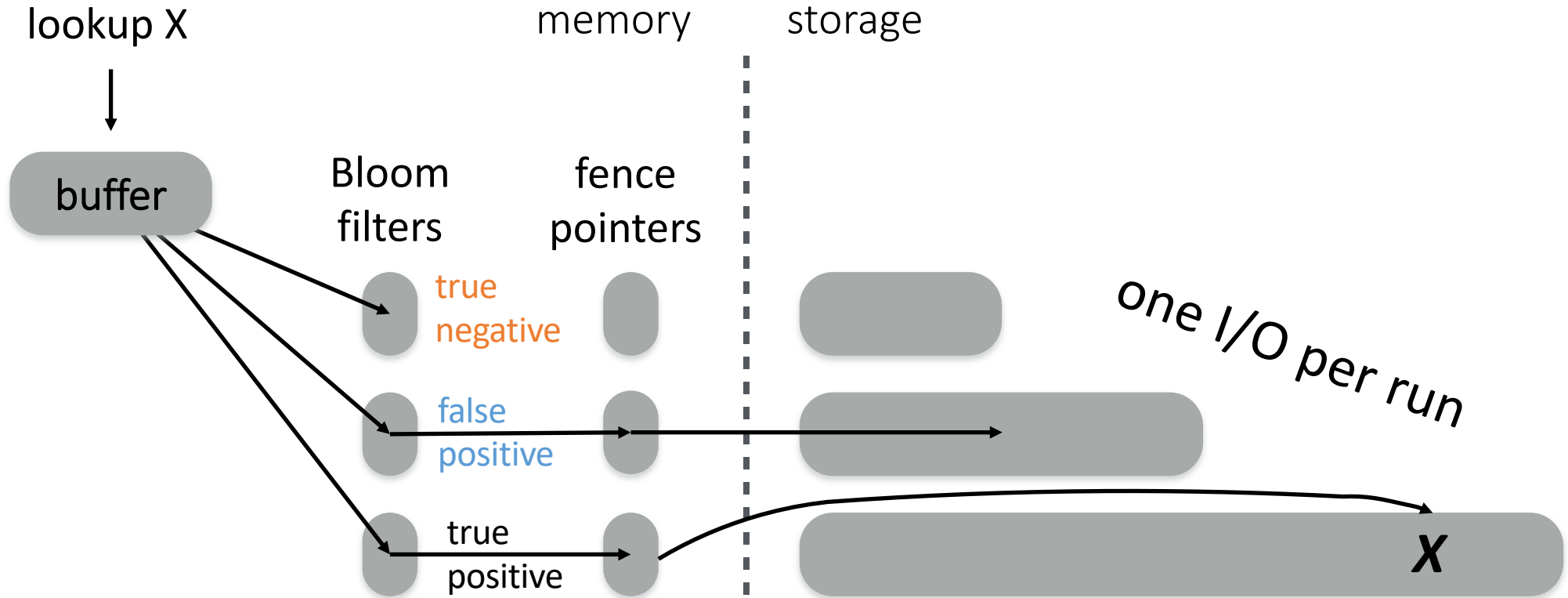
	Point Queries	Short Range Queries	Long Range Queries	Data Skew	Updates
B+ Trees	✓	✓	✓	✓	🌀
LSM Trees					
Radix Trees					
Hash Indexes					
Bitmap Indexes					
Scan Accelerators					

LSM-trees

Is LSM-tree an index at all?



LSM-trees



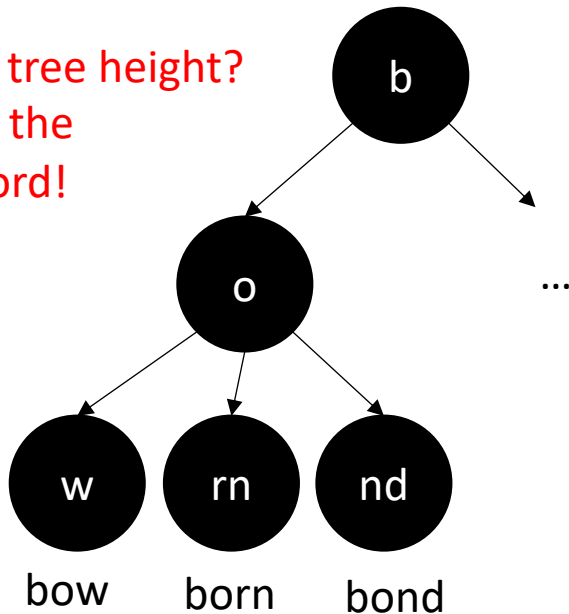
What are the possible index designs?

	Point Queries	Short Range Queries	Long Range Queries	Data Skew	Updates
B+ Trees	✓	✓	✓	✓	🌀
LSM Trees	✓	✗	🌀	✓	✓
Radix Trees					
Hash Indexes					
Bitmap Indexes					
Scan Accelerators					

Radix Trees (special case of tries and prefix B-Trees)

Idea: use common prefixes for internal nodes to reduce size/height!

Maximum tree height?
the size of the
longest word!



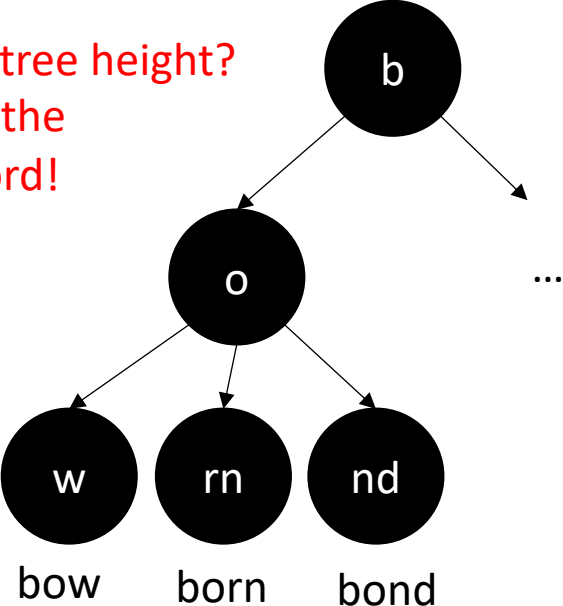
what about integer keys?

Radix Trees (special case of tries and prefix B-Trees)

Idea: use common prefixes for internal nodes to reduce size/height!

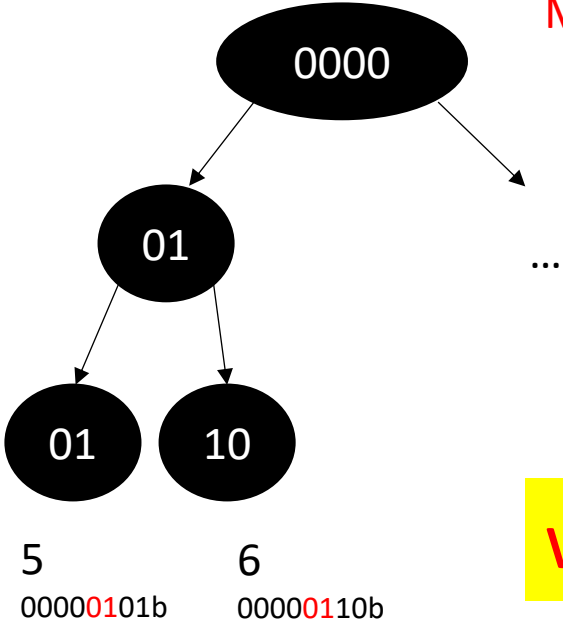
Binary representation of any domain can be used

Maximum tree height?
the size of the
longest word!



Maximum tree height?

8, that is, $\log_2(\max_domain_value)$
fixed worst case!



what about data skew?

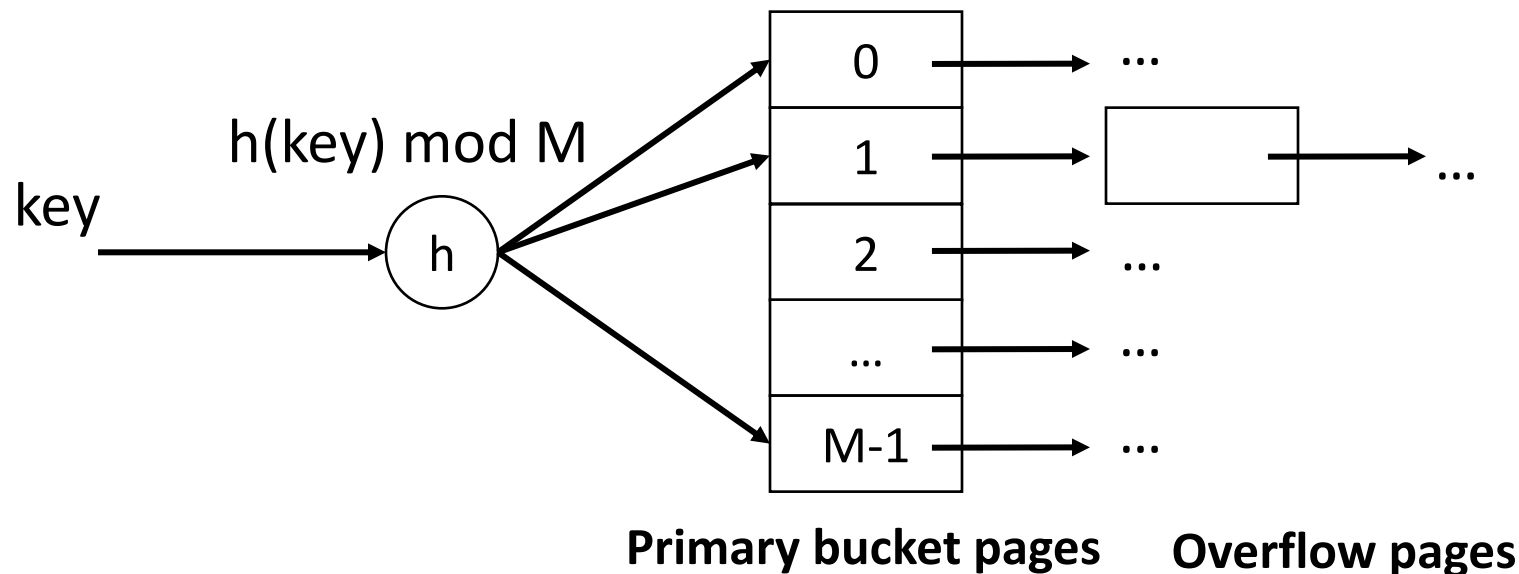
What are the possible index designs?

	Point Queries	Short Range Queries	Long Range Queries	Data Skew	Updates
B+ Trees	✓	✓	✓	✓	🌀
LSM Trees	✓	✗	🌀	✓	✓
Radix Trees	✓	✓	✓	✗	🌀
Hash Indexes					
Bitmap Indexes					
Scan Accelerators					

Hash Indexes (static hashing)

#primary bucket pages fixed, allocated sequentially, never de-allocated; overflow pages if needed

$h(k) \bmod M =$ bucket to insert data entry with key k (M : #buckets)

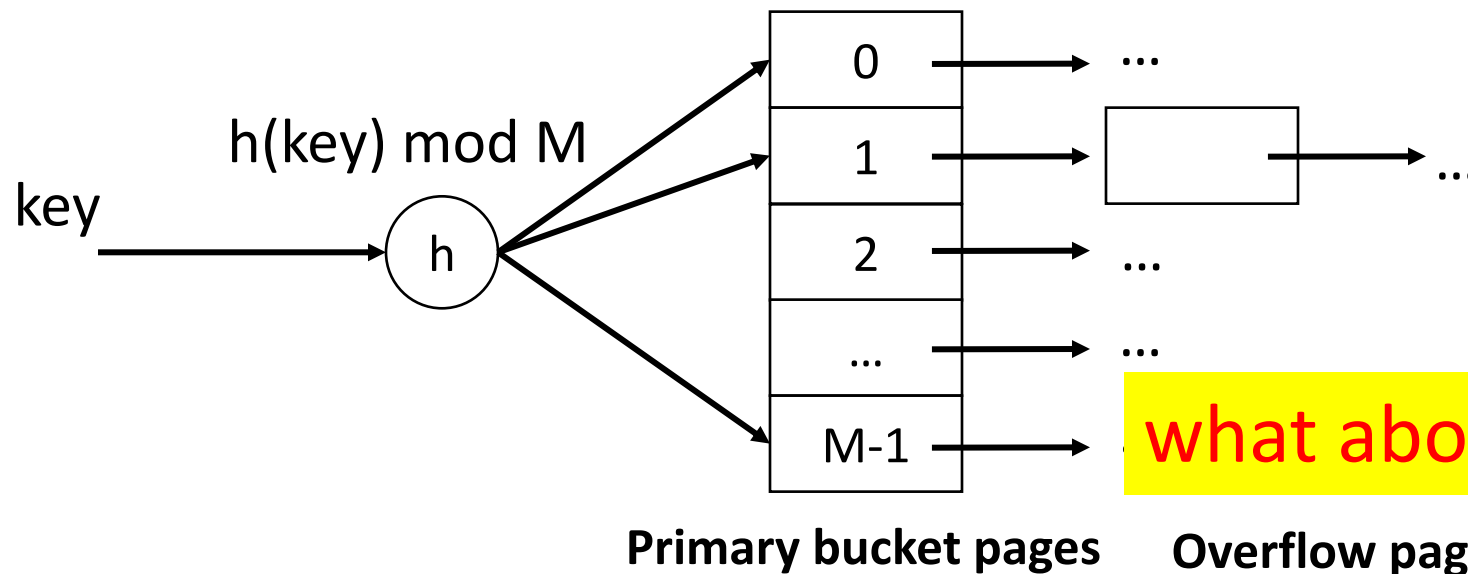


what if I have skew in the data set (or a bad hash function)?

Hash Indexes (static hashing)

#primary bucket pages fixed, allocated sequentially, never de-allocated; overflow pages if needed

$h(k) \bmod M =$ bucket to insert data entry with key k (M : #buckets)



what about range queries?

what if I have skew in the data set (or a bad hash function)?

What are the possible index designs?

	Point Queries	Short Range Queries	Long Range Queries	Data Skew	Updates
B+ Trees	✓	✓	✓	✓	🌀
LSM Trees	✓	✗	🌀	✓	✓
Radix Trees	✓	✓	✓	✗	🌀
Hash Indexes	✓	🌀	✗	✗	✓
Bitmap Indexes					
Scan Accelerators					

Bitmap Indexes

Column A	A=10	A=20	A=30
30	0	0	1
20	0	1	0
30	0	0	1
10	1	0	0
20	0	1	0
10	1	0	0
30	0	0	1
20	0	1	0

Speed & Size

- Compact representation of query result
- Query result is readily available

Bitvectors

- Can leverage fast Boolean operators
- Bitwise AND/OR/NOT faster than looping over meta data

Bitmap Indexes

Column A	A=10	A=20	A=30
30	0	0	1
20	0	1	0
30	0	0	1
10	1	0	0
20	0	1	0
10	1	0	0
30	0	0	1
20	0	1	0

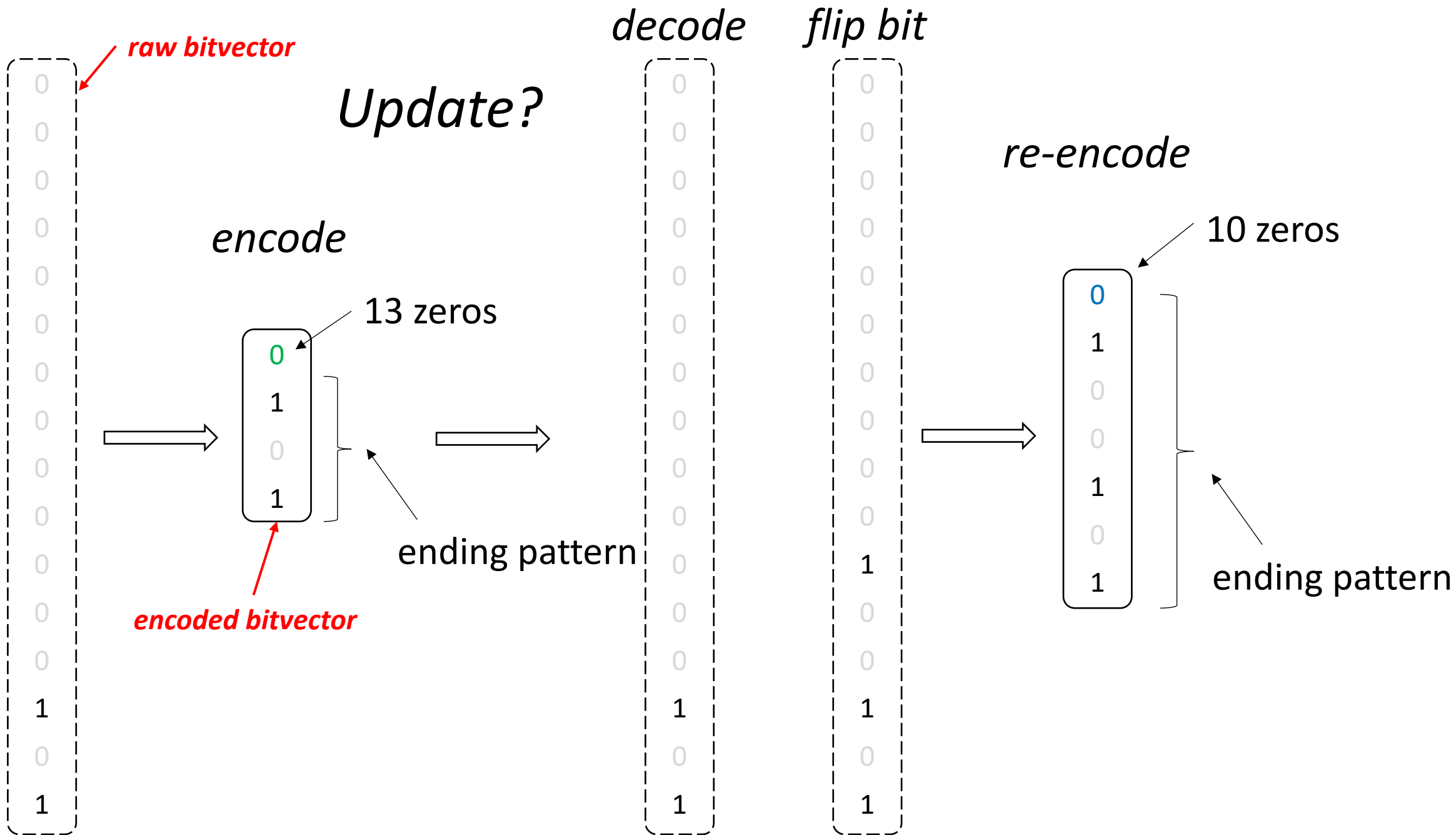
Index Size

- Space-inefficient for domains with large cardinality
- Addressed by bitvector encoding/compression

core idea: *run-length encoding* in prior work

encoded bitvectors

what about updates?



What are the possible index designs?

	Point Queries	Short Range Queries	Long Range Queries	Data Skew	Updates
B+ Trees	✓	✓	✓	✓	🌀
LSM Trees	✓	✗	🌀	✓	✓
Radix Trees	✓	✓	✓	✗	🌀
Hash Indexes	✓	🌀	✗	✗	✓
Bitmap Indexes	✓	🌀	🌀	🌀	✗
Scan Accelerators					

Scan Accelerators

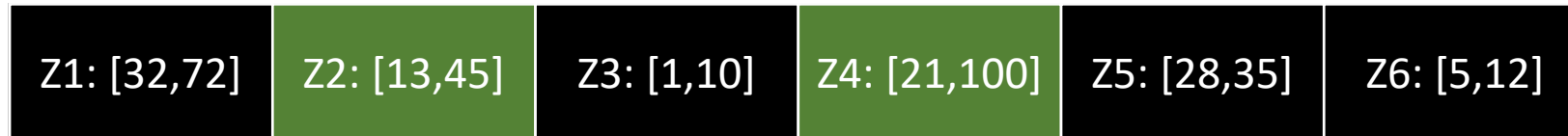
Zonemaps

Search for 25

Z1: [32,72]	Z2: [13,45]	Z3: [1,10]	Z4: [21,100]	Z5: [28,35]	Z6: [5,12]
-------------	-------------	------------	--------------	-------------	------------

Scan Accelerators

Zonemaps

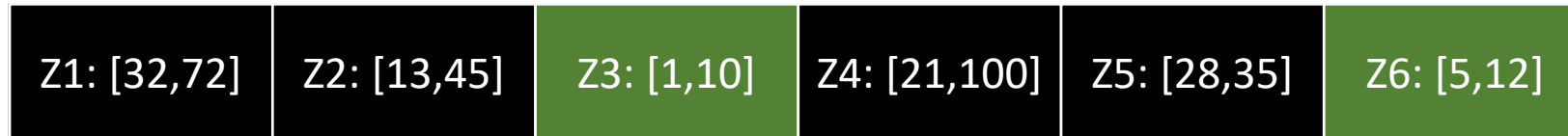


Search for 25

Search for [5,11]

Scan Accelerators

Zonemaps



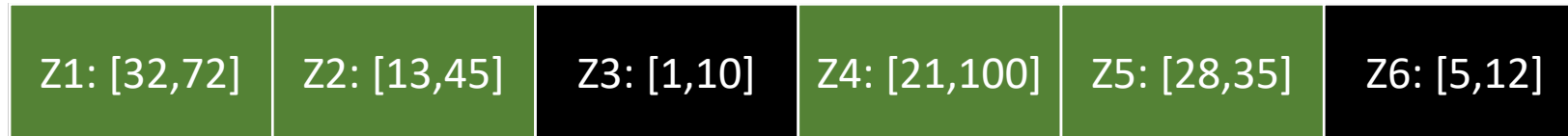
Search for 25

Search for [5,11]

Search for [31,46]

Scan Accelerators

Zonemaps



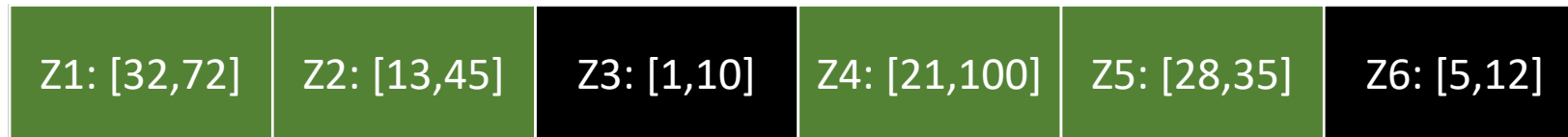
Search for 25

Search for [5,11]

Search for [31,46]

Scan Accelerators

Zonemaps



Search for 25

Search for [5,11]

Search for [31,46]

if data were sorted:



Search for 25

Search for [5,11]

Search for [31,46]

Scan Accelerators

Zonemaps

Z1: [32,72]	Z2: [13,45]	Z3: [1,10]	Z4: [21,100]	Z5: [28,35]	Z6: [5,12]
-------------	-------------	------------	--------------	-------------	------------

Search for 25

Search for [5,11]

Search for [31,46]

if data were sorted:

Z1: [1,15]	Z2: [16,30]	Z3: [31,50]	Z4: [50,67]	Z5: [68,85]	Z6: [85,100]
------------	-------------	-------------	-------------	-------------	--------------

Search for 25

Search for [5,11]

Search for [31,46]

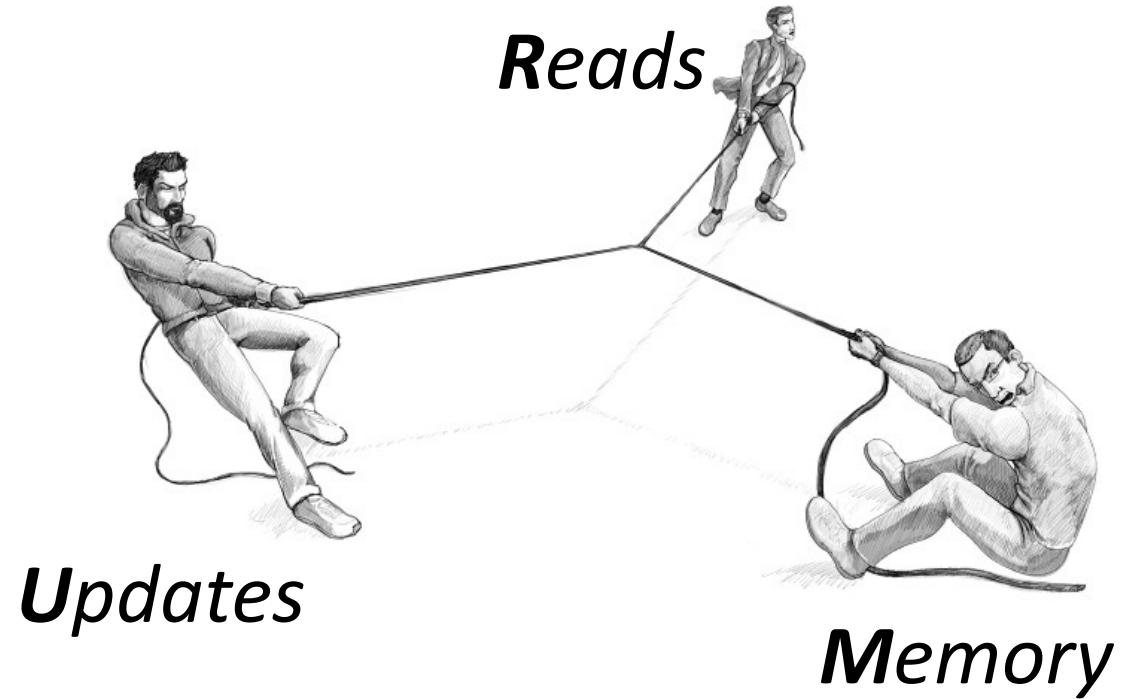
what if data is perfectly uniformly distributed?

Z1: [1,99]	Z2: [2,95]	Z3: [1,100]	Z4: [2,100]	Z5: [3,97]	Z6: [2,99]
------------	------------	-------------	-------------	------------	------------

What are the possible index designs?

	Point Queries	Short Range Queries	Long Range Queries	Data Skew	Updates
B+ Trees	✓	✓	✓	✓	⊘
LSM Trees	✓	✗	⊘	✓	✓
Radix Trees	✓	✓	✓	✗	⊘
Hash Indexes	✓	⊘	✗	✗	✓
Bitmap Indexes	✓	⊘	⊘	⊘	✗
Scan Accelerators	✗	⊘	✓	✓	⊘

data structure designs navigate a three-way tradeoff



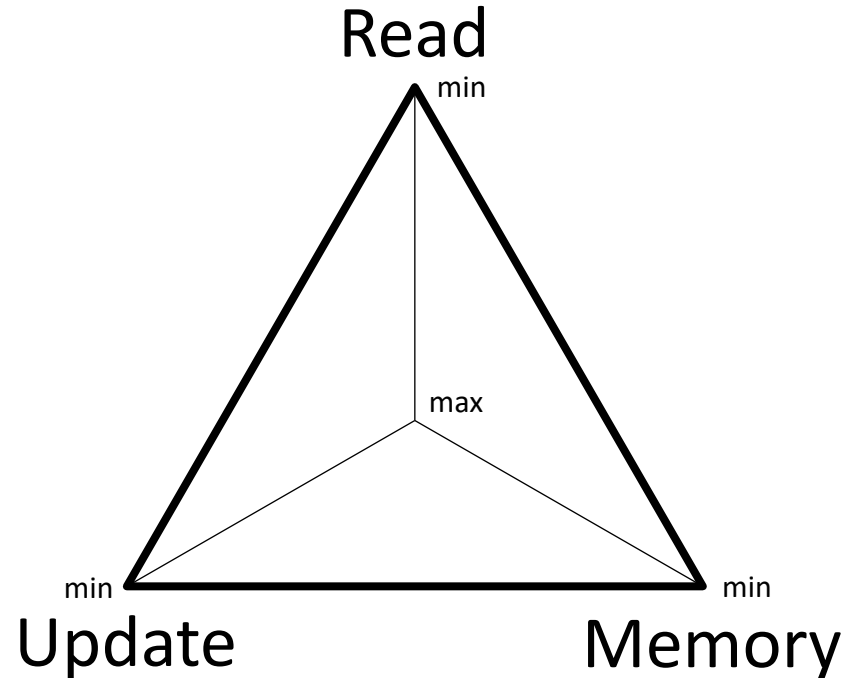
The RUM Conjecture

every access method has a (quantifiable)

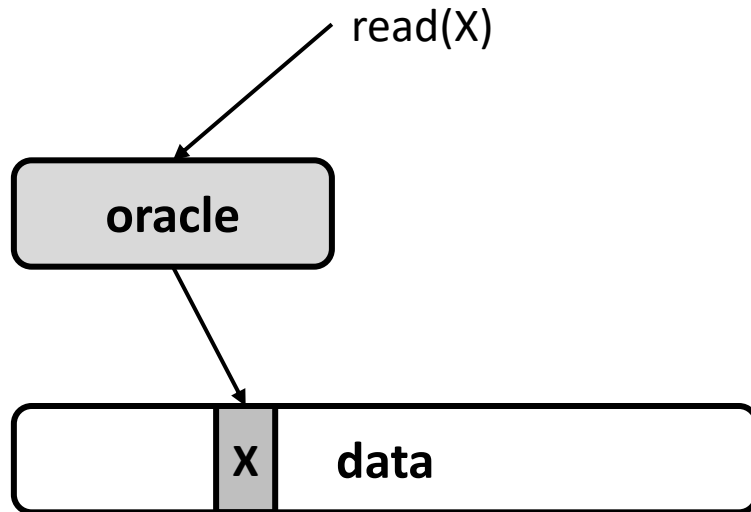
- read overhead
- update overhead
- memory overhead

the three of which form a competing triangle

we can optimize for two of the overheads at the expense of the third

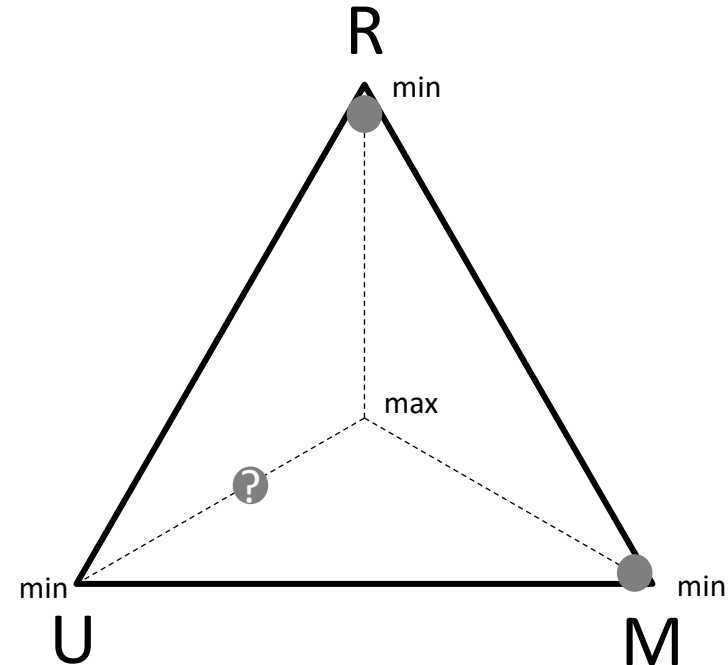


what would be an optimal read behavior?

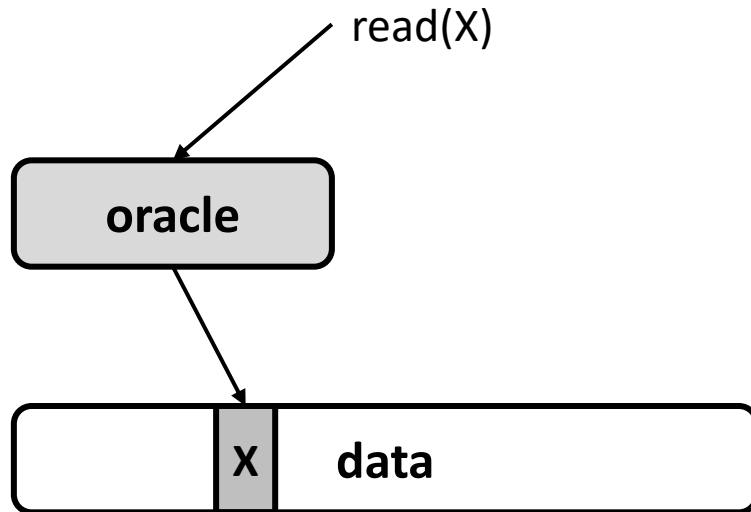


read(x) accesses only the bytes of object X

how *free* can an oracle be?

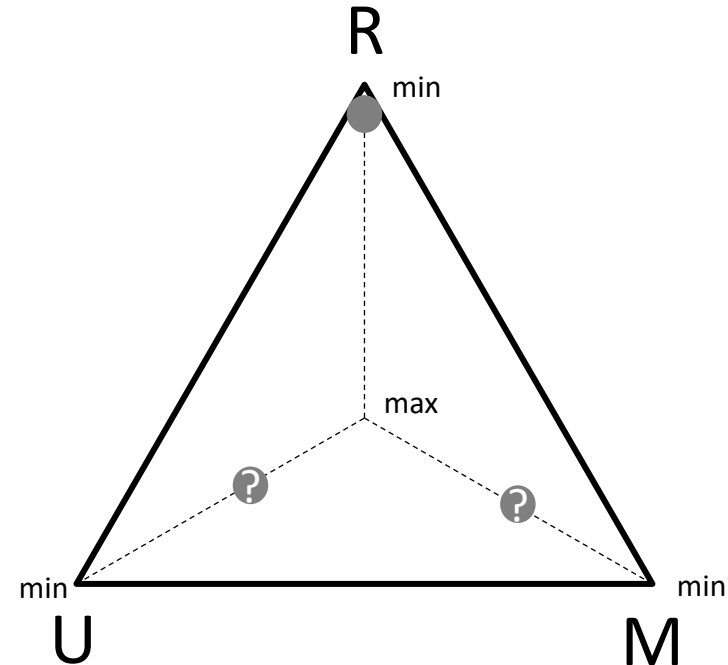


what would be an optimal read behavior?

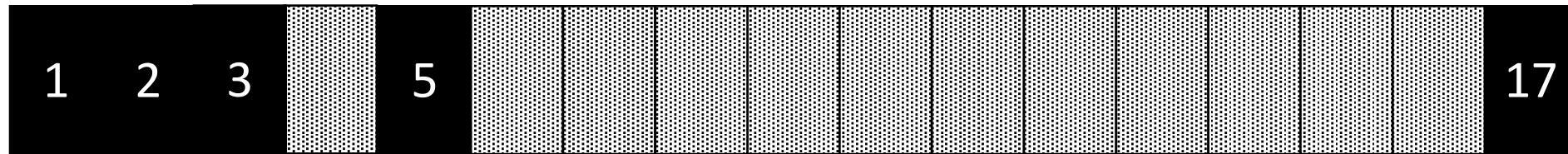


read(x) accesses only the bytes of object X

how *free* can an oracle be?



what would be an **optimal** read behavior?



update 17 -> 3

minimum read overhead

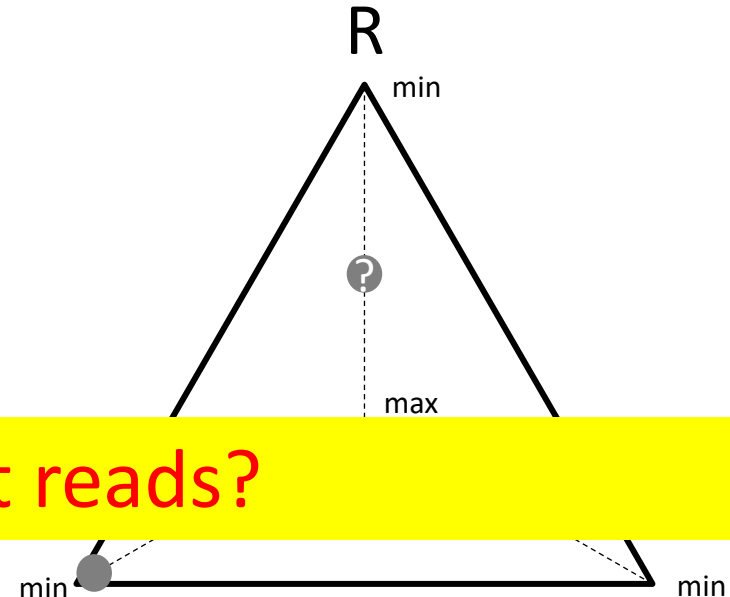
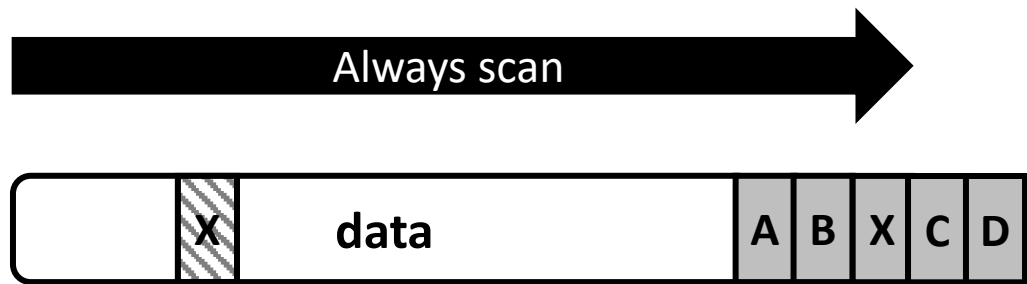
bound update overhead

unbounded memory overhead

what would be an optimal update behavior?

always *append*, and on update *invalidate*

update (X) changes the minimal number of bytes



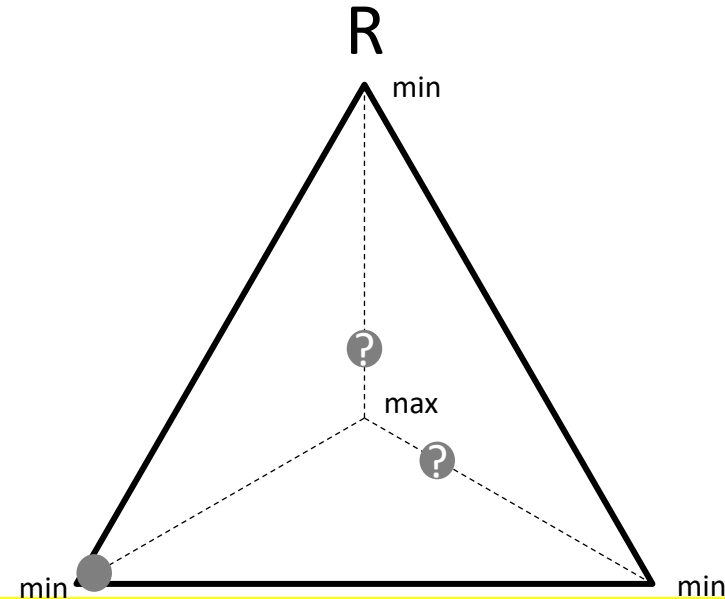
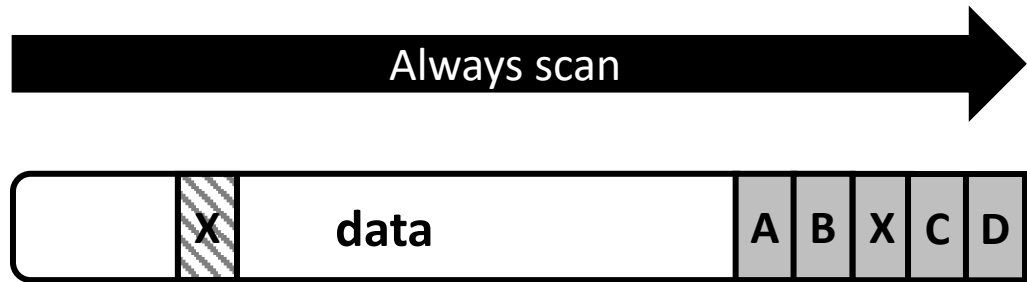
what about reads?

more data?

what would be an optimal update behavior?

always *append*, and *invalidate* on update

update (X) changes the minimal number of bytes

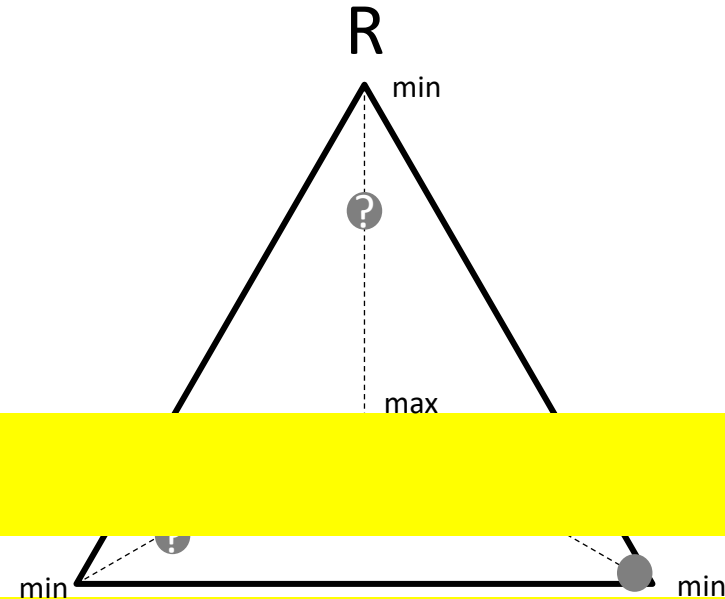
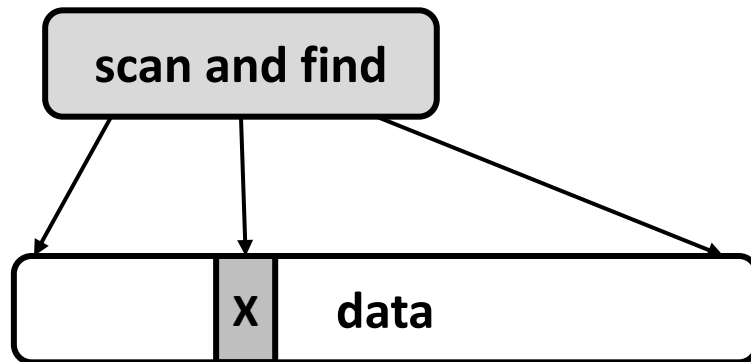


higher read and memory overhead

what would be an optimal memory overhead?

no metadata whatsoever, would result in the smallest memory footprint

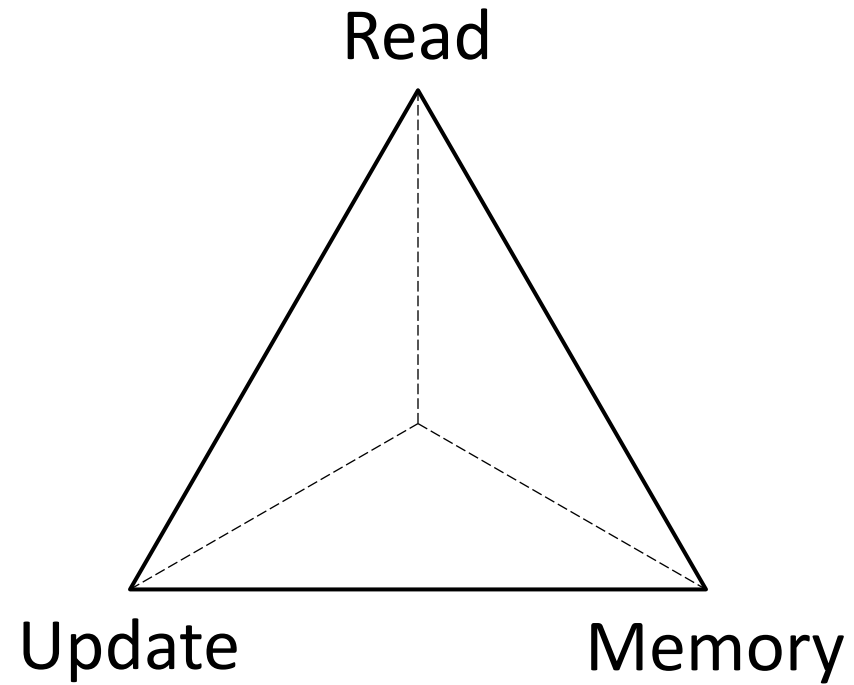
scan and in-place updates



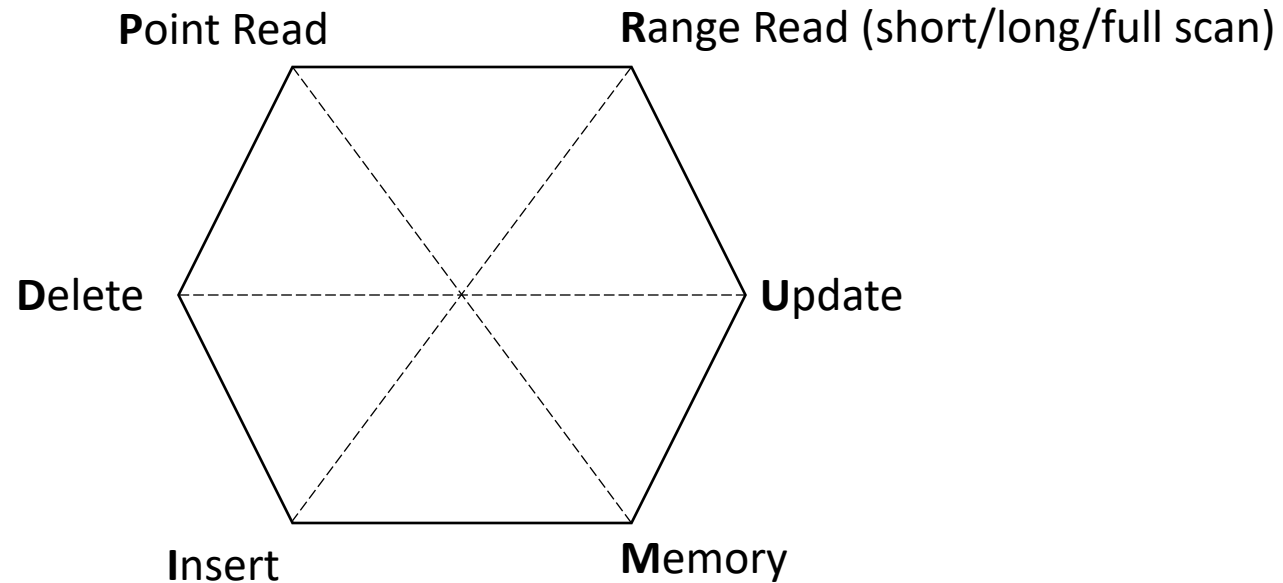
No!

do we need to reach the optimal(s)?

are there only three overheads?



are there only three overheads?

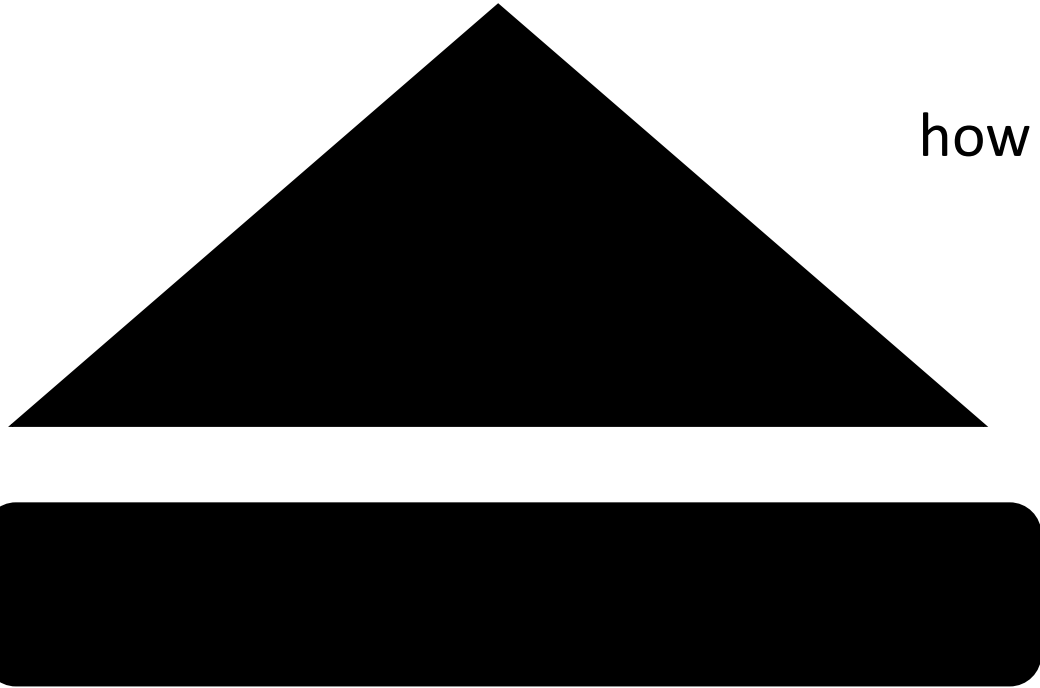


PyRUMID overheads

how to decide how to *design* a data structure?

break it down to *design dimensions*

how to break down the *design* in independent *dimensions*?



how to physically organize the data?

how to search through the data?

can I accelerate search through metadata?

multiple levels of nested organization?

how to update or add new data?

how to exploit additional memory/storage?

should the above decisions be applied eagerly or lazily?

how to break down the *design* in independent *dimensions*?



global data organization

how to search through the data?

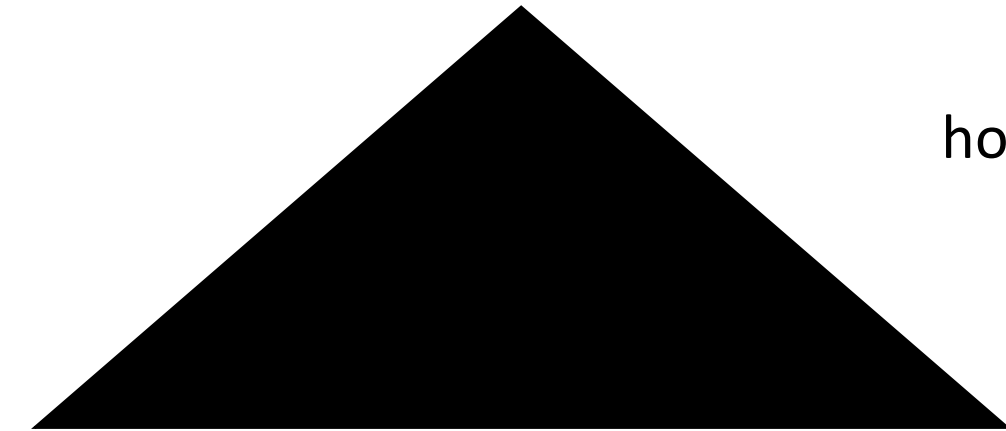
can I accelerate search through metadata?

multiple levels of nested organization?

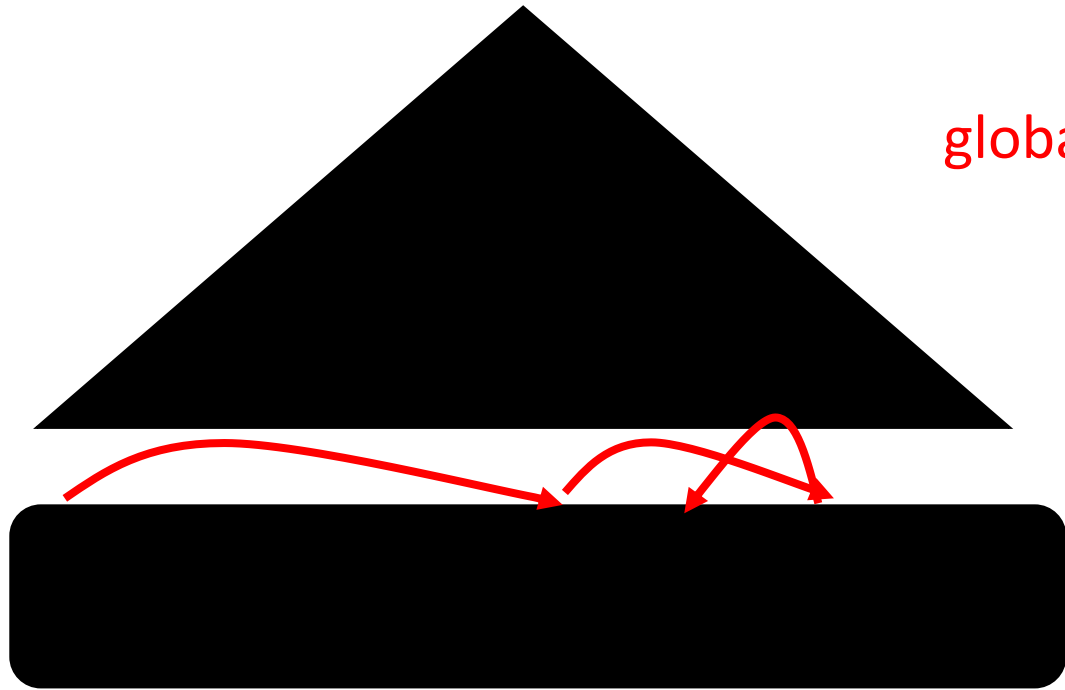
how to update or add new data?

how to exploit additional memory/storage?

should the above decisions be applied eagerly or lazily?



how to break down the *design* in independent *dimensions*?



global data organization

global search algorithm

can I accelerate search through metadata?

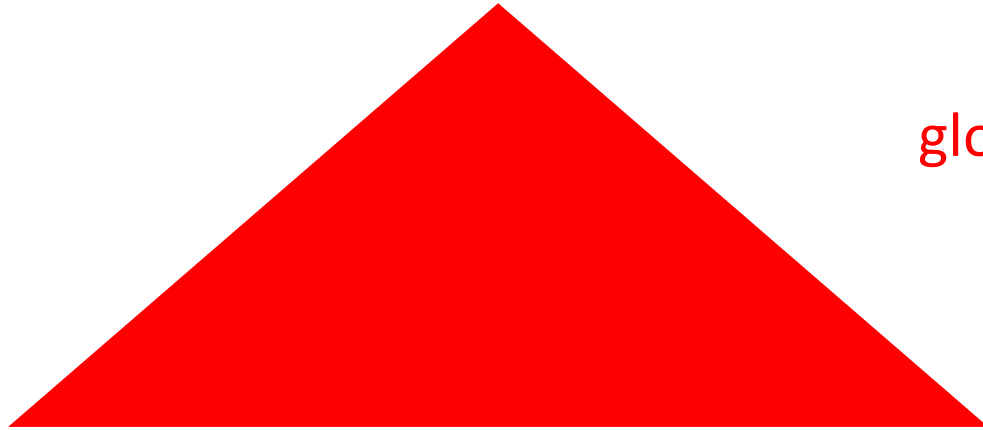
multiple levels of nested organization?

how to update or add new data?

how to exploit additional memory/storage?

should the above decisions be applied eagerly or lazily?

how to break down the *design* in independent *dimensions*?



global data organization

global search algorithm

metadata for searching

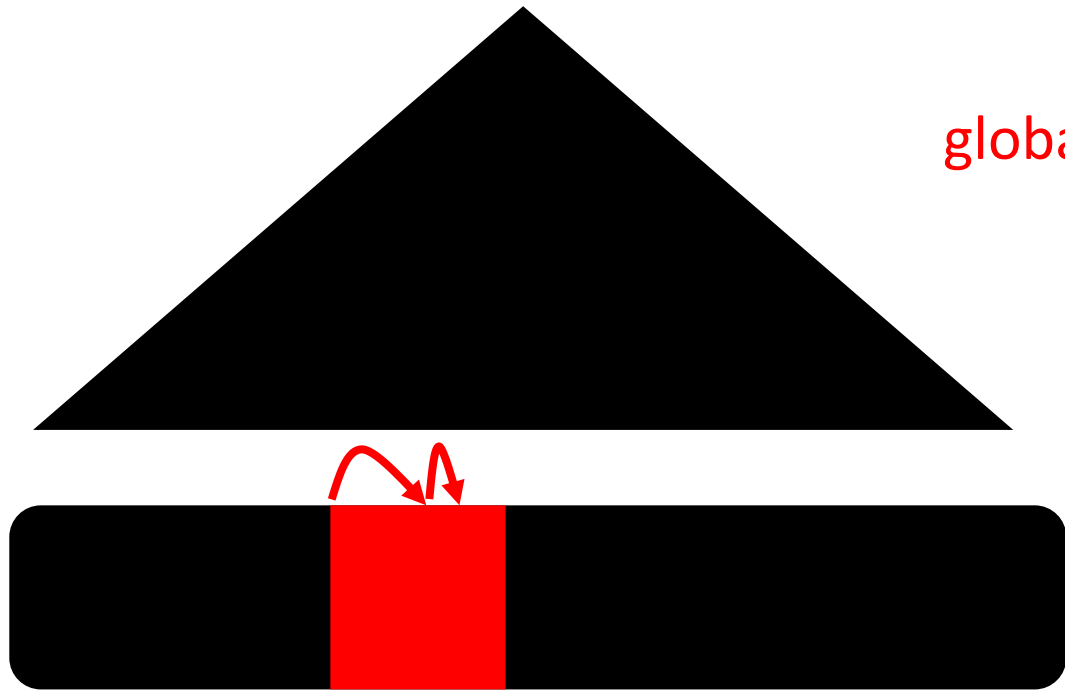
multiple levels of nested organization?

how to update or add new data?

how to exploit additional memory/storage?

should the above decisions be applied eagerly or lazily?

how to break down the *design* in independent *dimensions*?



global data organization

global search algorithm

metadata for searching

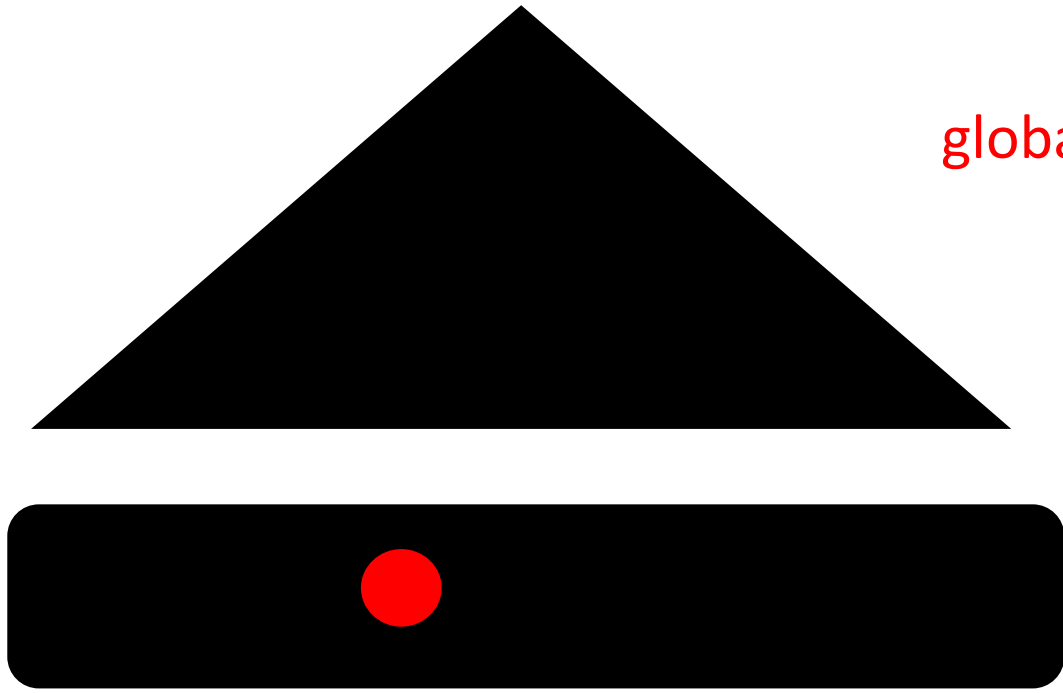
local data organization & search algorithm

how to update or add new data?

how to exploit additional memory/storage?

should the above decisions be applied eagerly or lazily?

how to break down the *design* in independent *dimensions*?



global data organization

global search algorithm

metadata for searching

local data organization & search algorithm

modification policy

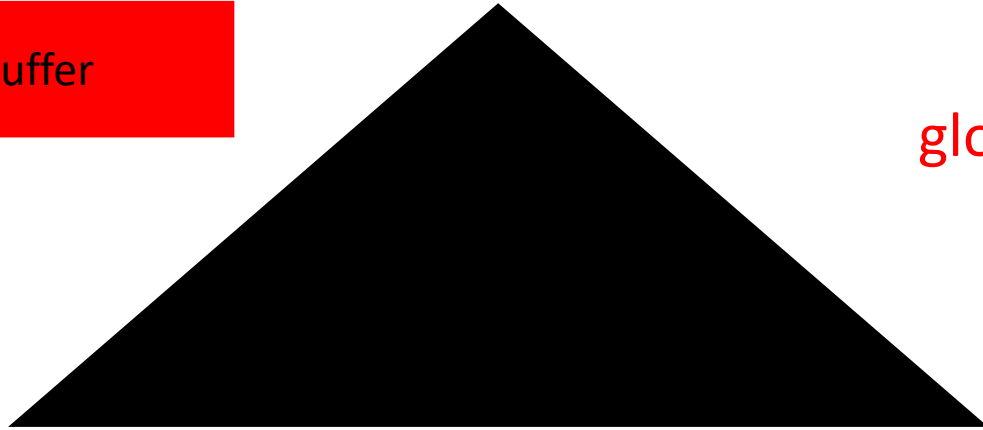
how to exploit additional memory/storage?

should the above decisions be applied eagerly or lazily?

how to break down the *design* in independent *dimensions*?



buffer



global data organization

global search algorithm

metadata for searching

local data organization & search algorithm

modification policy

batching via buffering

should the above decisions be applied eagerly or lazily?

how to break down the *design* in independent *dimensions*?



data structures *design dimensions and their values*

global data organization

global search algorithm

metadata for searching

local data organization & search algorithm

modification policy

batching via buffering

adaptivity

how to break down *popular designs*
to those design decisions?

b+ trees

global data organization

global searching (algorithm or index)

local data organization

local search algorithm

modification policy

Workload?



point and range queries and some modifications



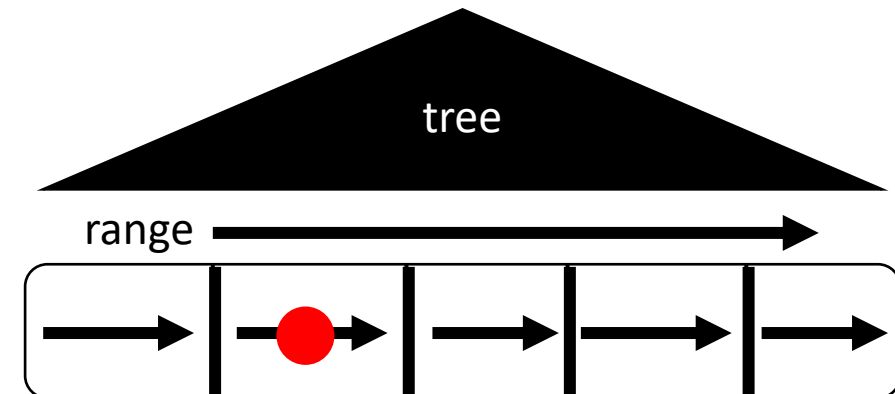
range partitioning

search tree

sorted

binary search / scan

in-place



insert optimized b+ trees

global data organization

global searching (algorithm or index)

local data organization

local search algorithm

modification policy

Workload?



increased number of modifications



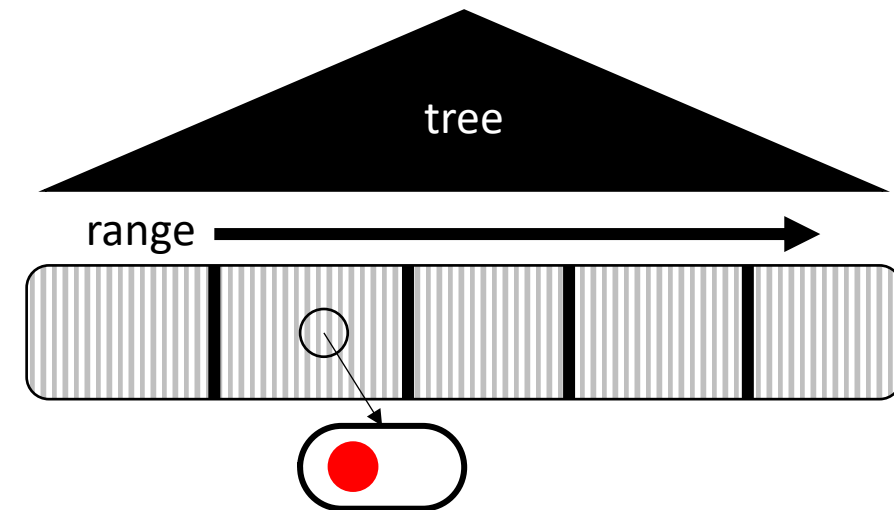
range partitioning

search tree

logging

binary search / scan

deferred in-place



static hashing

global data organization

global searching (algorithm or index)

local data organization

local search algorithm

modification policy

Workload?



point queries and modifications



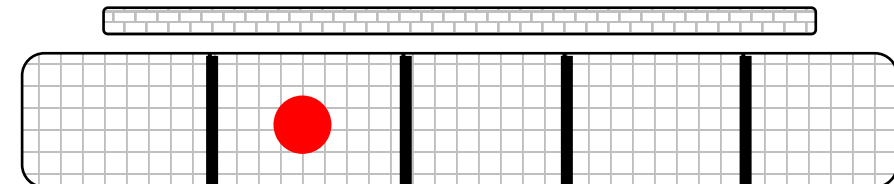
hash partitioning

direct addressing (hashing)

logging

scan

in-place



lsm-trees



global data organization

global searching (algorithm or index)

local data organization

local search algorithm

modification policy

Workload?



modification-heavy with point and range queries

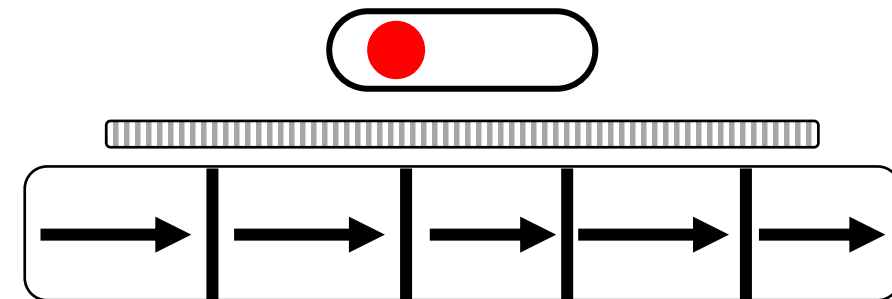
partitioned logging

filter indexing

sorted

binary / data-driven search

out-of-place



Projects

- A) LSM-tree implementation
- B) Bufferpool implementation
- C) WA quantification for LSMs on SSDs
- D) Range deletes in LSMs
- E) Query-driven compactions in LSMs
- F) Finding the optimal compaction strategy in LSMs to reduce WA
- G) Finding the optimal indexing granularity for LSMs
- H) Evaluating sorting algorithms for varied data sortedness
- I) Measuring robustness in SplinterDB
- J) Boosting join performance in Postgres for skewed correlation
- K) Benchmarking large-scale graph processing systems

CS 561: Data Systems Architectures

Introduction to Indexing:

Trees, Tries, Hashing, Bitmaps: The Whole Design Space

Dr. Subhadeep Sarkar

<https://bu-disc.github.io/CS561/>