



NYU

**TANDON SCHOOL
OF ENGINEERING**

Introduction to Index Structures

CS6086: Principles of Database Systems
NYU Tandon



NEW YORK UNIVERSITY



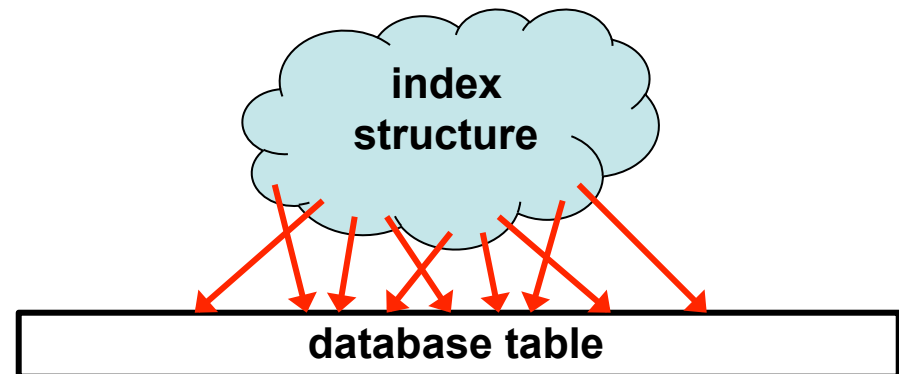
NYU

TANDON SCHOOL
OF ENGINEERING

■ Indexing:

- An index is a data structure that supports efficient retrieval of “certain types of tuples” from a table

- Setup:



- Without index:

- Need to scan entire table
- Or binary search if sorted



NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

■ Problem:

“Find records of a certain type quickly”

- **Equality Queries**

$\sigma_{ssn = 619441789}(\text{Employees})$

- **Range Queries**

$\sigma_{10 \leq age \leq 20}(\text{Employees})$

These are the most interesting types of conditions

Plus AND and OR of several such conditions

Goal: faster processing of SQL queries (selects and joins)



NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

■ Various types of indexes

1. Clustered (Primary) or Nonclustered (Secondary)
2. Sparse or Dense
3. Single Key or Composite or Multi-Dimensional

Different implementation (algorithms)

- Tree Structures (ordered)
- Hashing (non-ordered)
- Others



NEW YORK UNIVERSITY

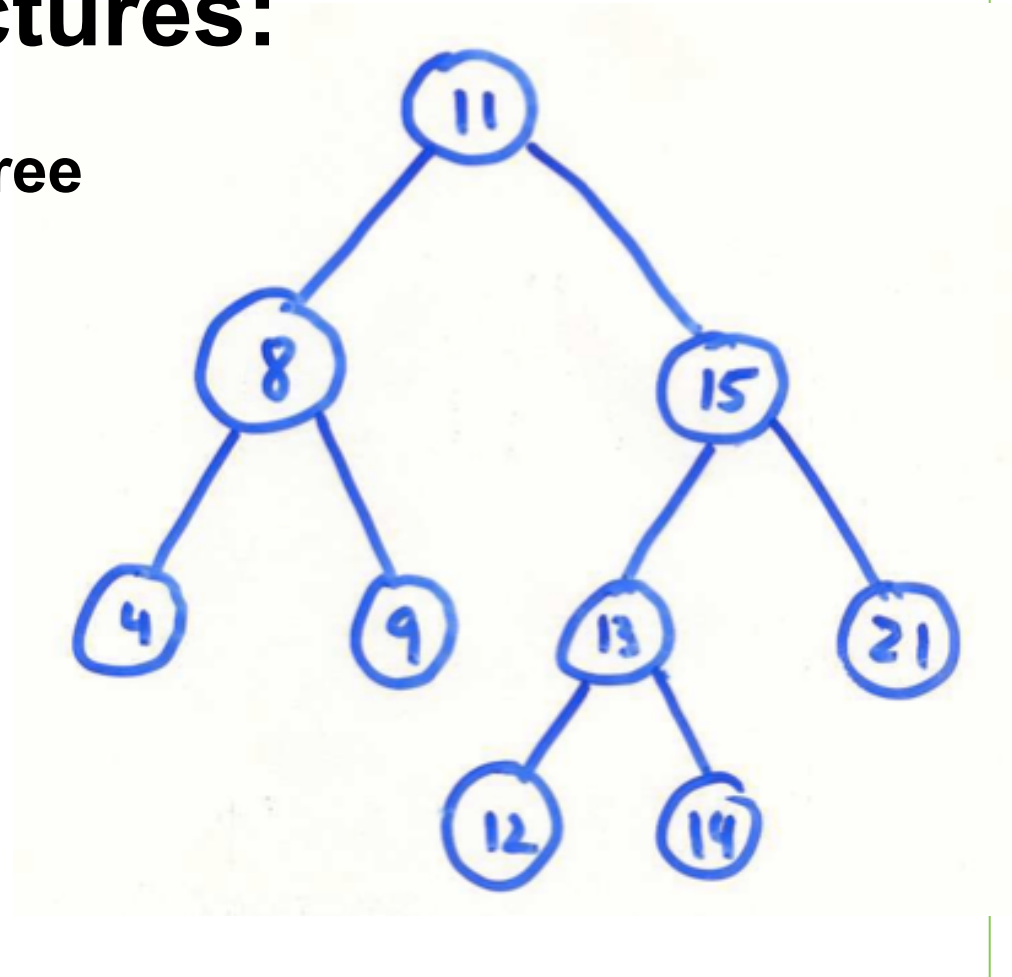


NYU

TANDON SCHOOL
OF ENGINEERING

■ Indexing Structures:

E.g., Binary Search Tree



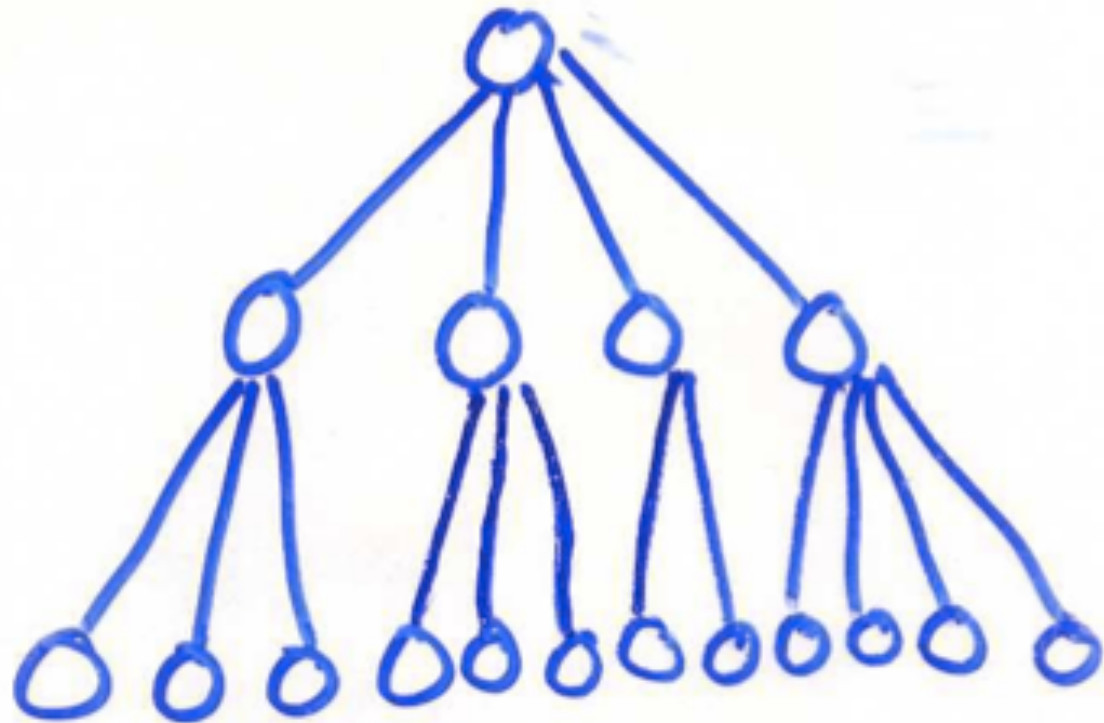
NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

B⁺-Tree:



- Degree $\gg 2$ (why?)
- Data in leaves only (why?)
- Same depth everywhere (how?)



NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

■ Compare: Trees with N Elems in Leaves

- Binary Tree with N keys:
 - $\log_2 N$ levels
 - Each internal node:
 - 20 bytes? (2 pointers + 1 key)
 - $N = 10^6 \Rightarrow 20$ levels
- Degree-100 tree:
 - $\log_{100} N$ levels
 - Each internal node:
 - 1200 (1196) bytes (100 pointers + 99 keys)
 - $N = 10^6 \Rightarrow 3$ levels
- \Rightarrow This is much better in the disk case!



NEW YORK UNIVERSITY

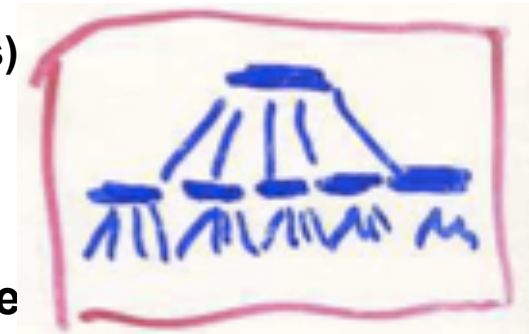
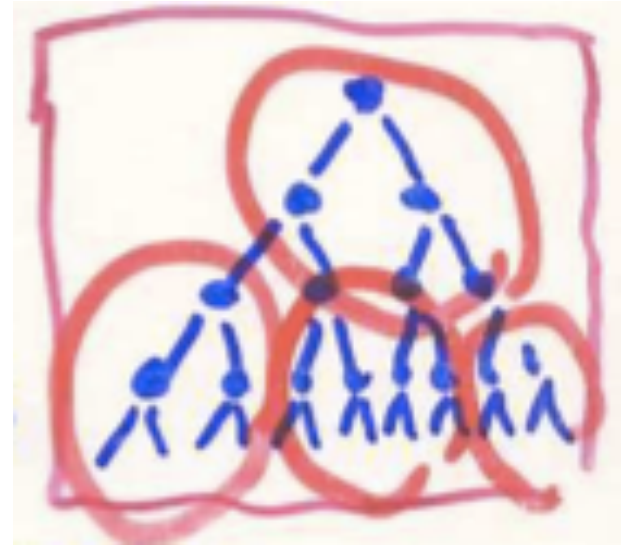


NYU

TANDON SCHOOL
OF ENGINEERING

Compare: Trees with N Elems in Leaves

- Binary Tree with N keys:
 - $\log_2 N$ levels
 - Each internal node:
 - 20 bytes? (2 pointers + 1 key)
 - $N = 10^6 \Rightarrow 20$ levels
- Degree-100 tree:
 - $\log_{100} N$ levels
 - Each internal node:
 - 1200 (1196) bytes (100 pointers + 99 keys)
 - $N = 10^6 \Rightarrow 3$ levels
- \Rightarrow This is much better in the disk case!
- .. like collapsing several binary levels into one



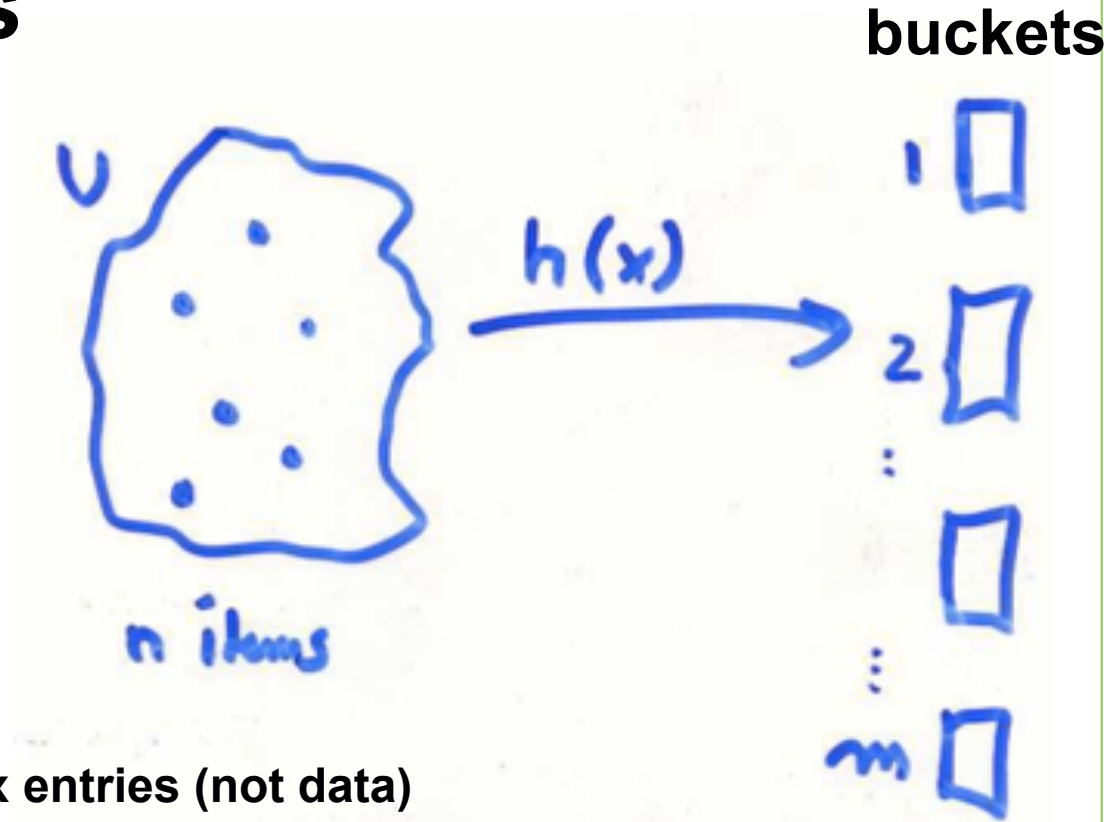
NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

Hash Tables



- x = key value
- $h: U \rightarrow [1,2,3,\dots,m]$
- buckets contain index entries (not data)
- No range queries supported
- Exception: small enumerable ranges



NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

■ Index entries

- An index entry is a pair where:

KEY	POINTER
-----	---------

- Key is a key value
 - Pointer is a record ID (usually not a physical pointer)
 - Every object in the DB has an RID.
- Index Construction:
 - Make an index entry for each tuple in the table
 - Insert index entries into index structure (B⁺-tree, hash, etc)



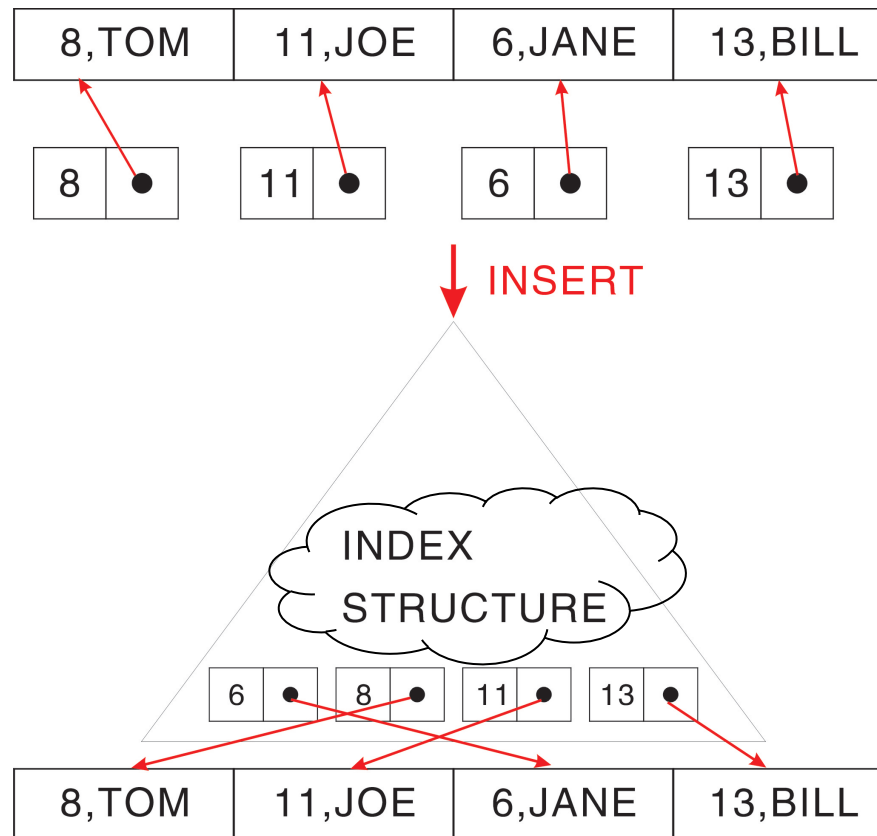
NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

Example: Schema is (age, name) with index on age



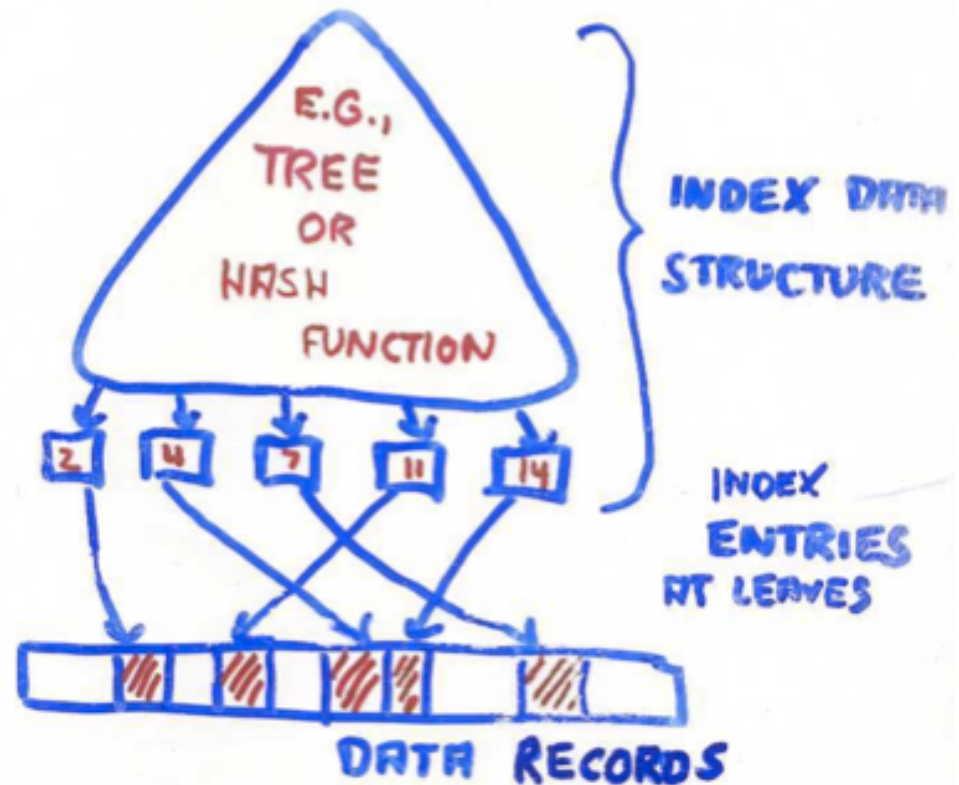
NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

Indexing: Basic Setup



Hides:

- Type of structure used (tree or hash)
- File/record organization used (data is not moved during indexing)



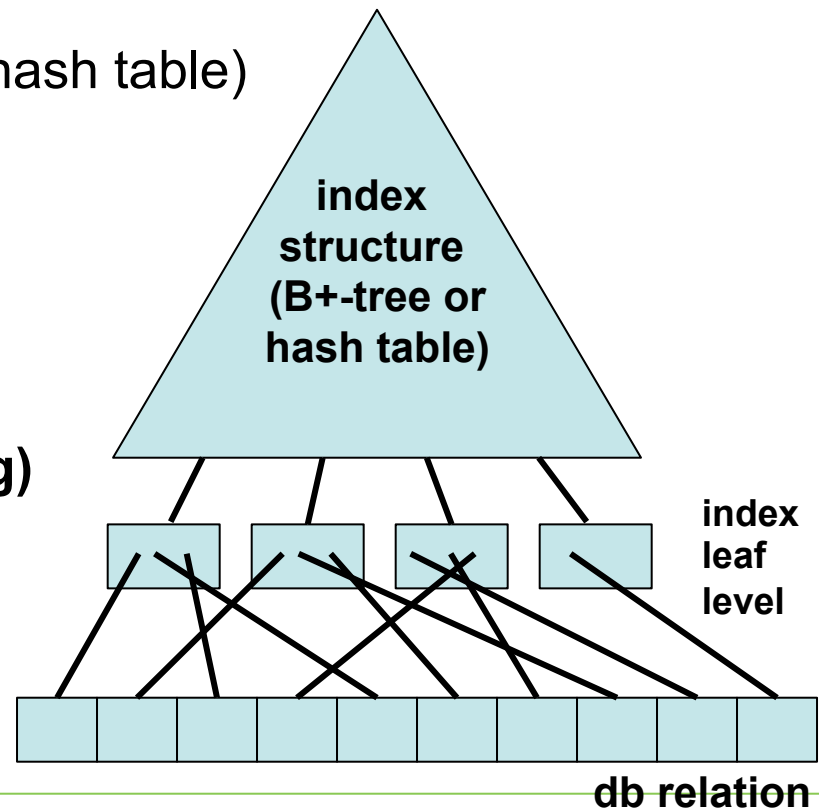
NEW YORK UNIVERSITY



Indexing: Basic Setup

Hides: index structure (B+-tree or hash table)

- Type of structure used (tree or hash)
- File/record organization (data not moved by indexing)

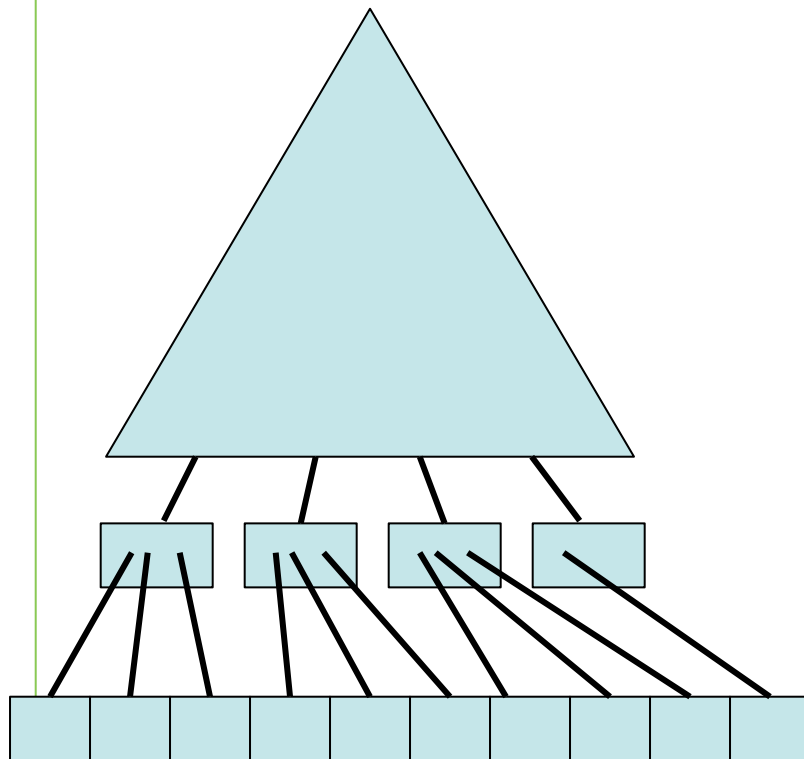




NYU

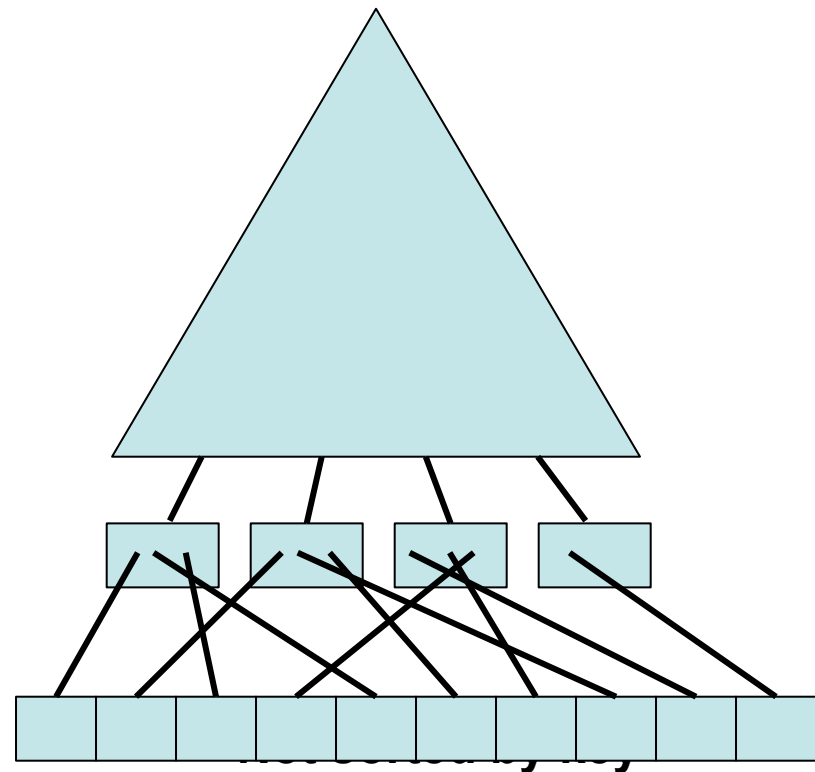
TANDON SCHOOL
OF ENGINEERING

Primary Index



relation sorted by key

Secondary Index



not sorted by key



NEW YORK UNIVERSITY



■ Primary Indexes

- Can have only one primary index per table
- Not (always) same as primary key! (e.g., SSN)
- Advantages
 - Efficient range search
 - Smaller index size possible (sparse index)
- Sparse Index:
 - Idea: don't create index entry for each row in table
 - E.g., only pointers to row at start of a data blocks





NYU

TANDON SCHOOL
OF ENGINEERING

■ Why Primary Indexes?

- Table is already sorted by attribute
- Why not just do binary search?



NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

■ Why Primary Indexes?

- Table is already sorted by attribute
- Why not just do binary search?
- Answer: binary search is not I/O-efficient
- Table with 10^6 records: probe $\lg_2(10^6) \sim 20$ records
- B+-tree: fetch only 4-5 index nodes
- Also, top levels of tree are cached
- Binary search:

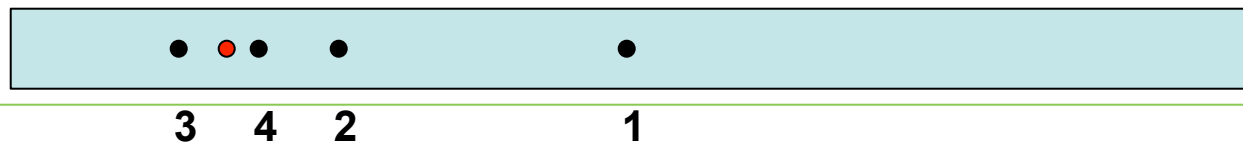


NEW YORK UNIVERSITY



■ Why Primary Indexes?

- Table is already sorted by attribute
- Why not just do binary search?
- Answer: binary search is not I/O-efficient
- Table with 10^6 records: probe $\lg_2(10^6) \sim 20$ records
- B+-tree: fetch only 4-5 index nodes
- Also, top levels of tree are cached
- Binary search:





■ Why Primary Indexes?

- Table is already sorted by attribute
- Why not just do binary search?
- Answer: binary search is not I/O-efficient
- Table with 10^6 records: probe $\lg_2(10^6) \sim 20$ records
- B+-tree: fetch only 4-5 index nodes
- Also, top levels of tree are cached
- Binary search: cache commonly probed items?



1/4 1/8 1/2

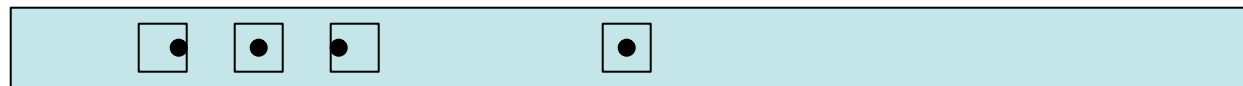
1





Why Primary Indexes?

- Table is already sorted by attribute
- Why not just do binary search?
- Answer: binary search is not I/O-efficient
- Table with 10^6 records: probe $\lg_2(10^6) \sim 20$ records
- B+-tree: fetch only 4-5 index nodes
- Also, top levels of tree are cached
- Binary search: caching is inefficient



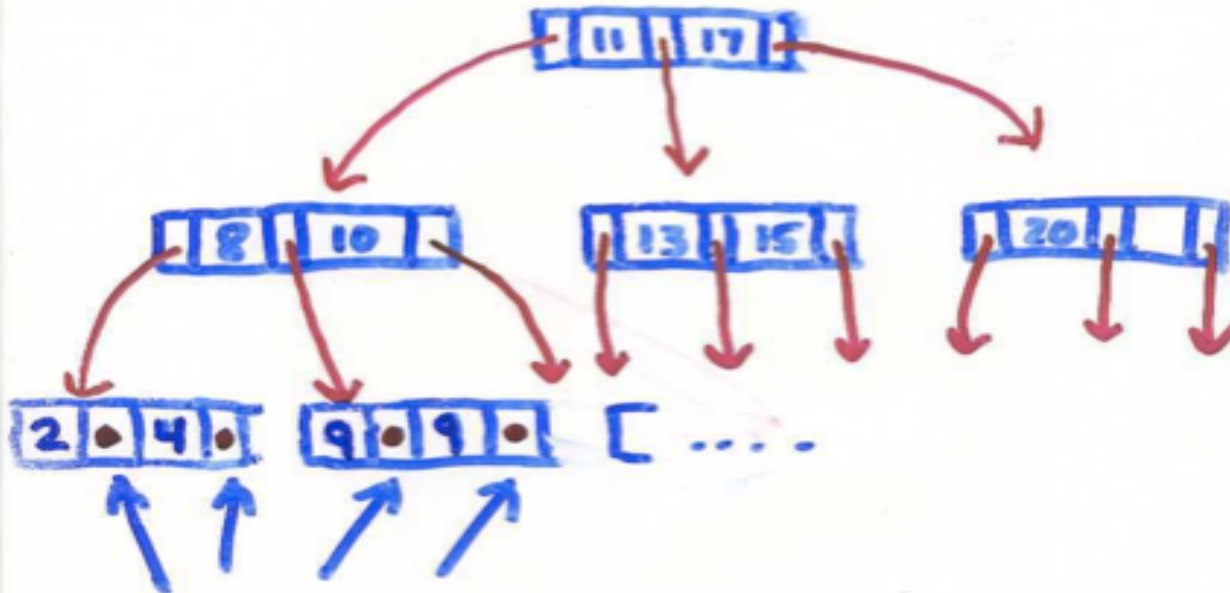


NYU

TANDON SCHOOL
OF ENGINEERING

Example: B⁺ - Tree

- “Switch to algorithms mode”



“PAYLOAD” (Pointer to record) in index entries



NEW YORK UNIVERSITY



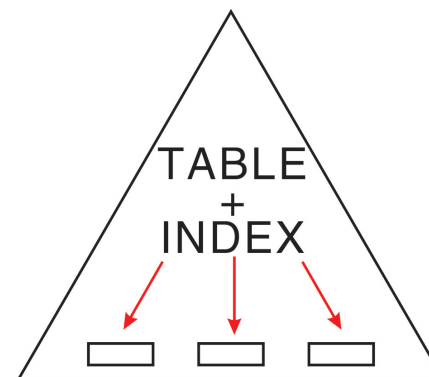
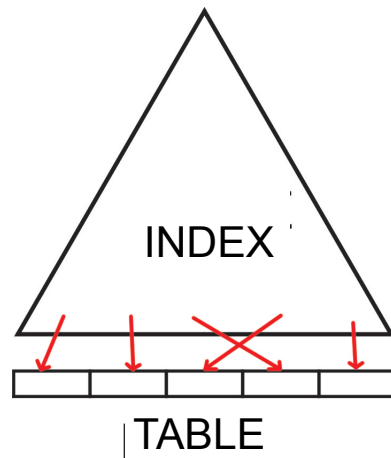
NYU

TANDON SCHOOL
OF ENGINEERING

Types of Data Structures Used

- Tree – Based
- Hash – Based
- Other (often multi-dim)

Index vs. File Organization:



➔ We assume index organization



NEW YORK UNIVERSITY



Types of indexes

- **Primary / Secondary & Clustered / Unclustered**
 - primary (clustered) means table sorted on index key
 - only one primary index per table
- **Dense / Sparse**
 - sparse only possible if clustered index
- **Single Attribute / Composite**
 - single attribute
 - attribute 1, attribute 2, ...
 - True multi-attribute (multi-dimensional) index

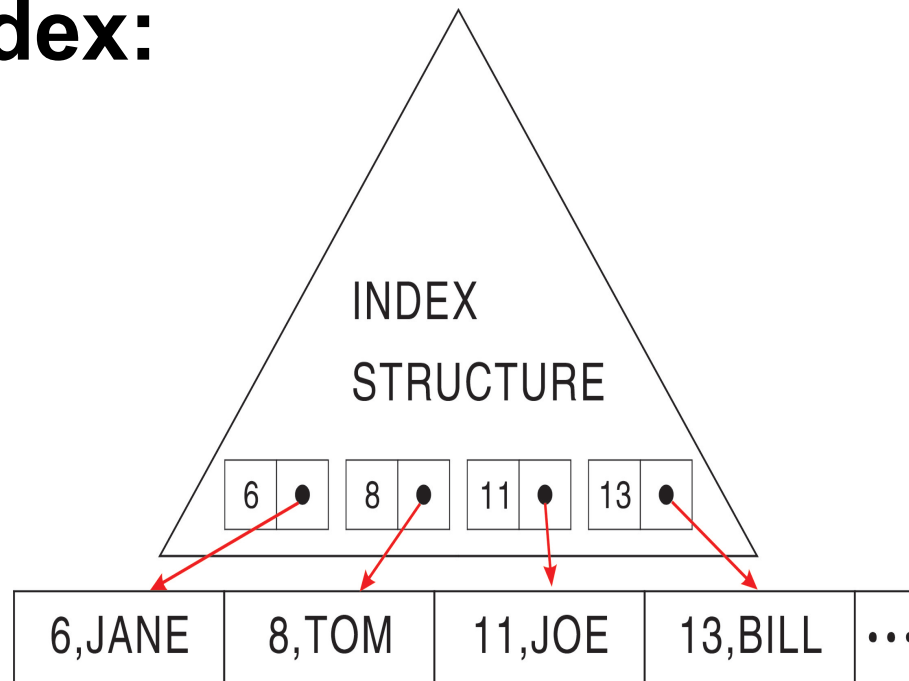




NYU

TANDON SCHOOL
OF ENGINEERING

Clustered Index:



- No crossing pointers
- Do we really need index ?
- Do we need all index entries?

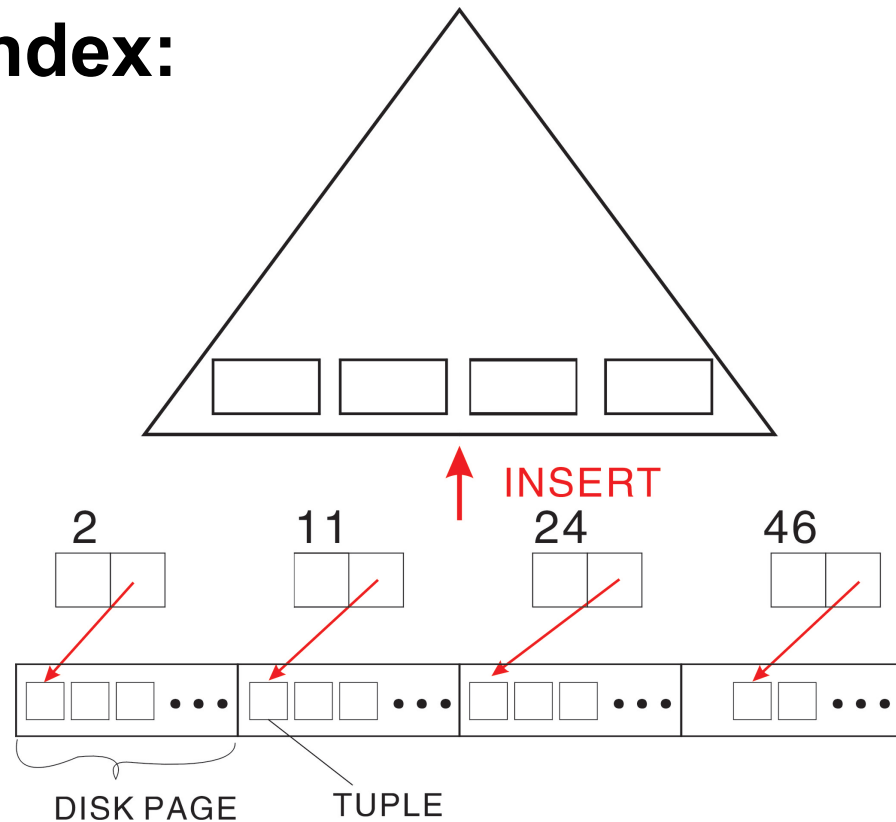


NEW YORK UNIVERSITY



Sparse Clustered Index:

- One entry per page
- Smaller Index
- Fewer levels
- But need to search inside pages





NYU

TANDON SCHOOL
OF ENGINEERING

■ B⁺-Trees in a Nutshell:

- Each node \approx disk block
- Typically 4 KB
- Internal nodes
 - $\lceil n/2 \rceil \leq c \leq n$ # of child pointers
 - $\lceil n/2 \rceil - 1 \leq k \leq n-1$ # of key values
- Root node: like internal, but at least 2 kids & 1 key
- Leaf nodes
 - $\lceil (n-1)/2 \rceil \leq k \leq n-1$
of key values and also # of record pointers

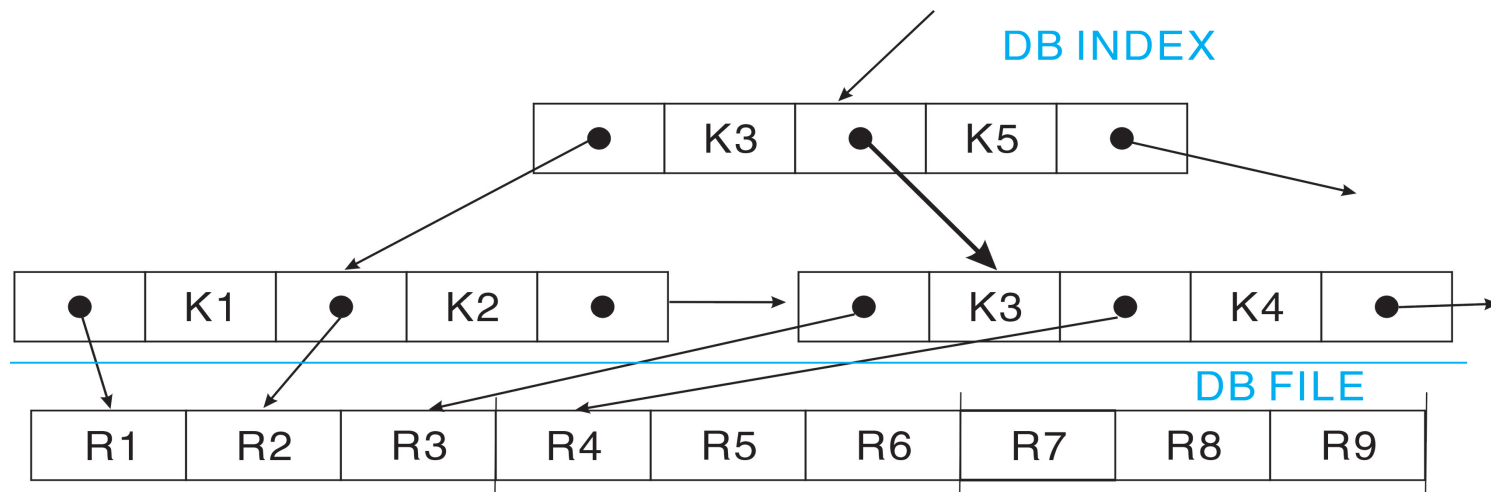


NEW YORK UNIVERSITY



Example:

- 4 KB node size
- 8 – byte integer keys
- 8 – byte pointers
- Choose $n = 250$ so that $250 \cdot 8 + 249 \cdot 8 = 3992$ bytes ≈ 4 KB (including some overhead)





NYU

TANDON SCHOOL
OF ENGINEERING

■ Insertion & Deletion:

- **Insert:**
 - If space left in leaf, just insert and done
 - Else, try to rebalance with sibling neighbor (optional)
 - If impossible, split node, insert new node as child in parent
 - If parent has no space, rebalance/split at parent level, etc.
- **Delete**
 - If leaf still $\geq \lceil n/2 \rceil$ children, just delete and done
 - Else, try to rebalance with sibling neighbor
 - If impossible, merge with sibling, delete node from parent
- **Rebalance:** balance data with sibling so each has same num



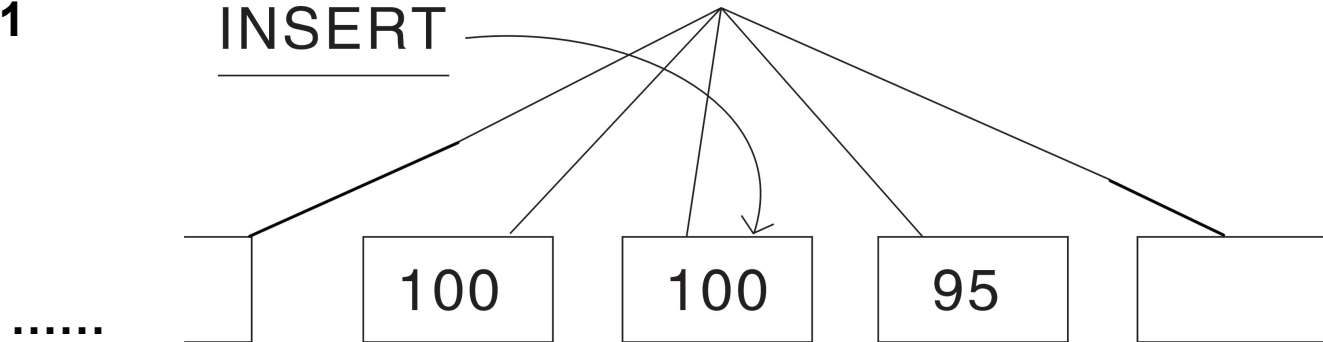
NEW YORK UNIVERSITY



- **Note: numbers are # of index entries in leaves**

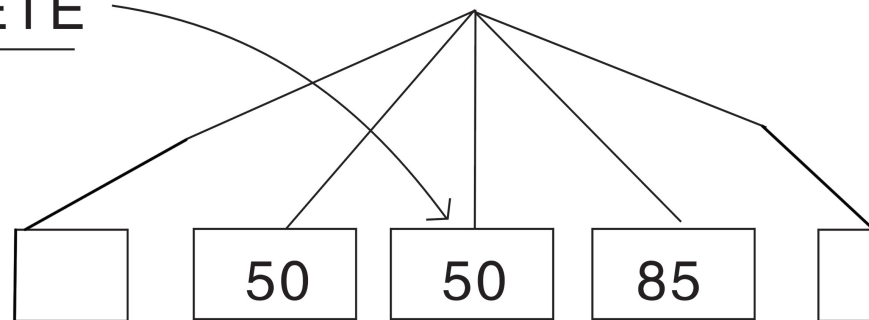
- $n = 101$

INSERT



- **Could balance with right neighbor (98 each afterwards) or not**
- **Delete: try to balance**

DELETE



- **Never go beyond immediate neighbor!**



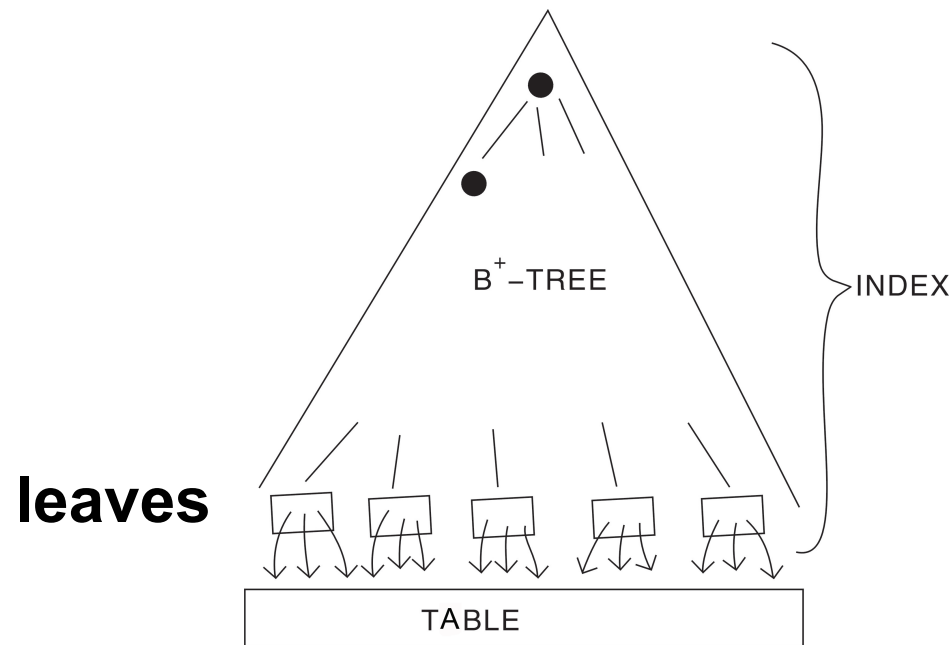


NYU

TANDON SCHOOL
OF ENGINEERING

File vs. Index Organization

- Example of B⁺ - Tree (Clustered) Index Organization



- Index and relation are separate!



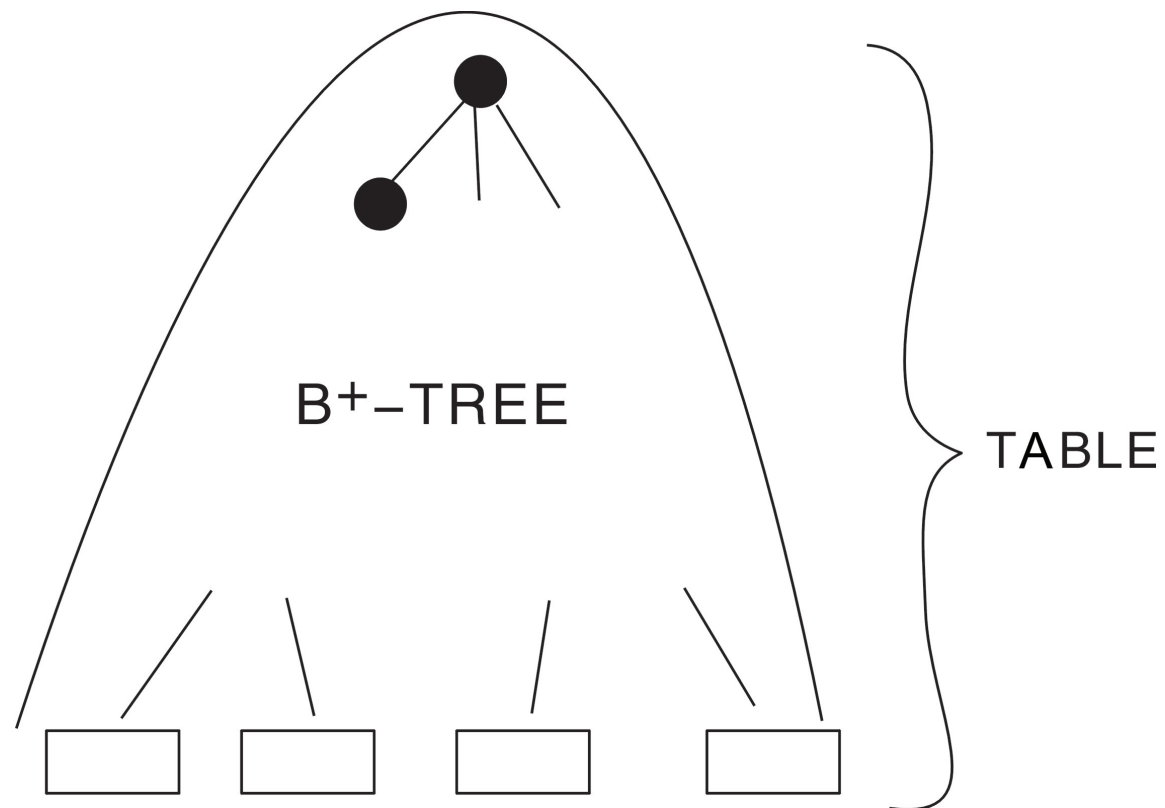
NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

B⁺-Tree File Organization



NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

■ B⁺-Tree File Organization

- **Tuples are contained in leaves (instead of pointers to tuples)**
- **Maybe one less indirection (1 fewer disk access to find tuple)**
- **But more leaf nodes (height may change)**
- **Can still add secondary indexes**
- **Similar for hash tables**
- **“The file is the Index”**

We mostly assume the other case: tuples not at leaf level



NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

Dense Non-Clustered B+-Tree Index:

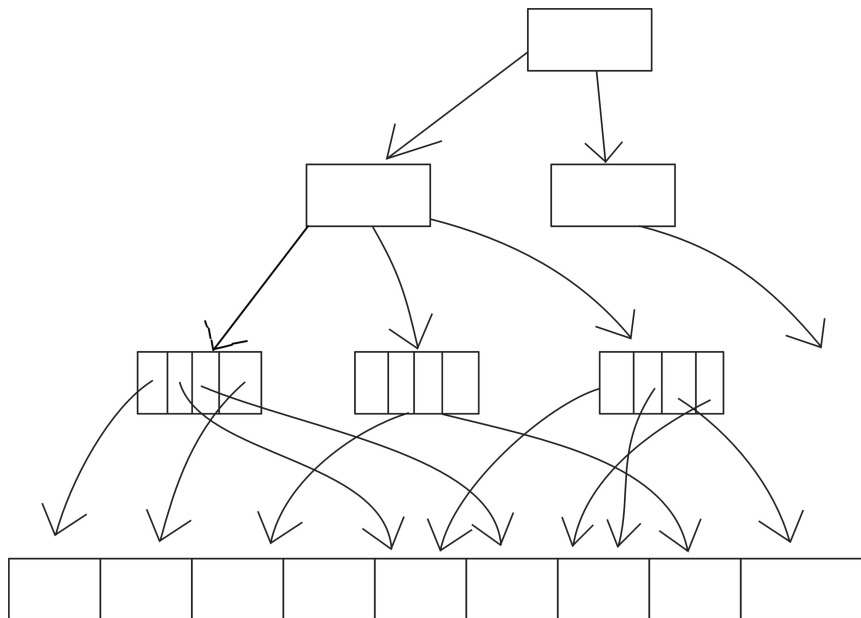
n = 250

1 root

16 ~ 64 nodes

4000 ~ 8000 nodes

1000000 records



- Leaf level: 1000000 (key, pointer) pairs
- Next level: 4000 ~ 8000 pairs



NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

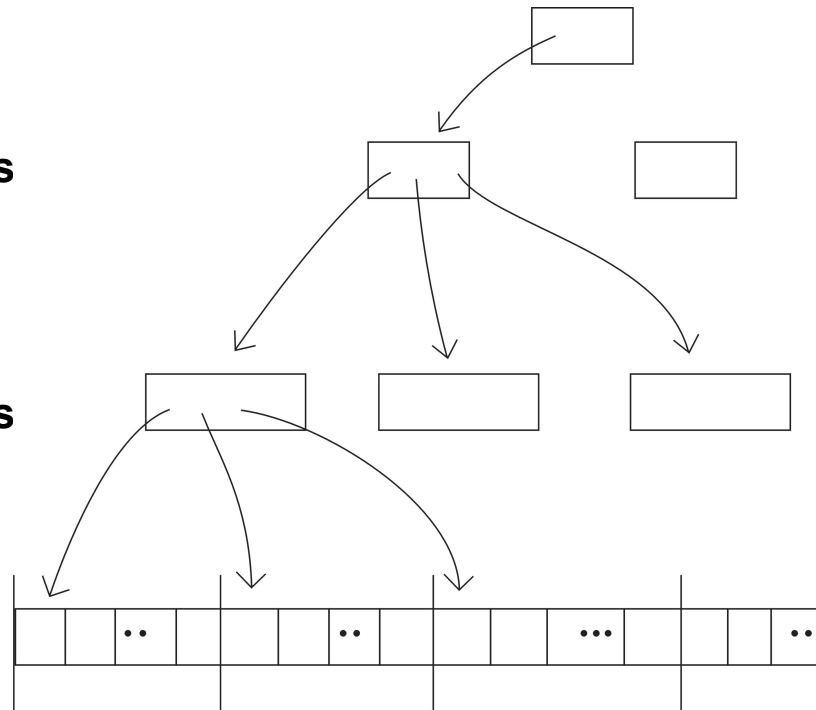
Sparse, Clustered B⁺-Tree Index:

- $n = 250$, 20 records/block

1 ~ 4 nodes

200 ~ 400 nodes

1000000 records in 50000 blocks



NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

■ Analysis of B⁺ - Tree Costs

- Assume occupancy factor
- E.G. 80% → degree 200 for $n = 250$
- Or upper / lower bounds as shown
- Also occupancy factor for DB file if sparse index
- Then count cost in terms of disk accesses and/or transfer costs



NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

■ ISAM vs. B+-Trees

- Index sequential access method
- Not as good for updates
- Slightly simpler
- But B+-trees overall better
- Conceptually: ISAM is “Indexing an index” but recursively
- We focus on B⁺-trees



NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

■ Example

- Disk with $T_s + T_r = 10\text{ms}$ and $tr = 10 \text{ MB/S}$
- Database with 1000000 records of size 100 bytes
- Assume block size 4 KB
 - ➔ 25000 blocks with 40 records each (100 % occupancy)
and $T_B = 10\text{ms} + 0.4\text{ms} = 10.4\text{ms}$ (time per table)



NEW YORK UNIVERSITY

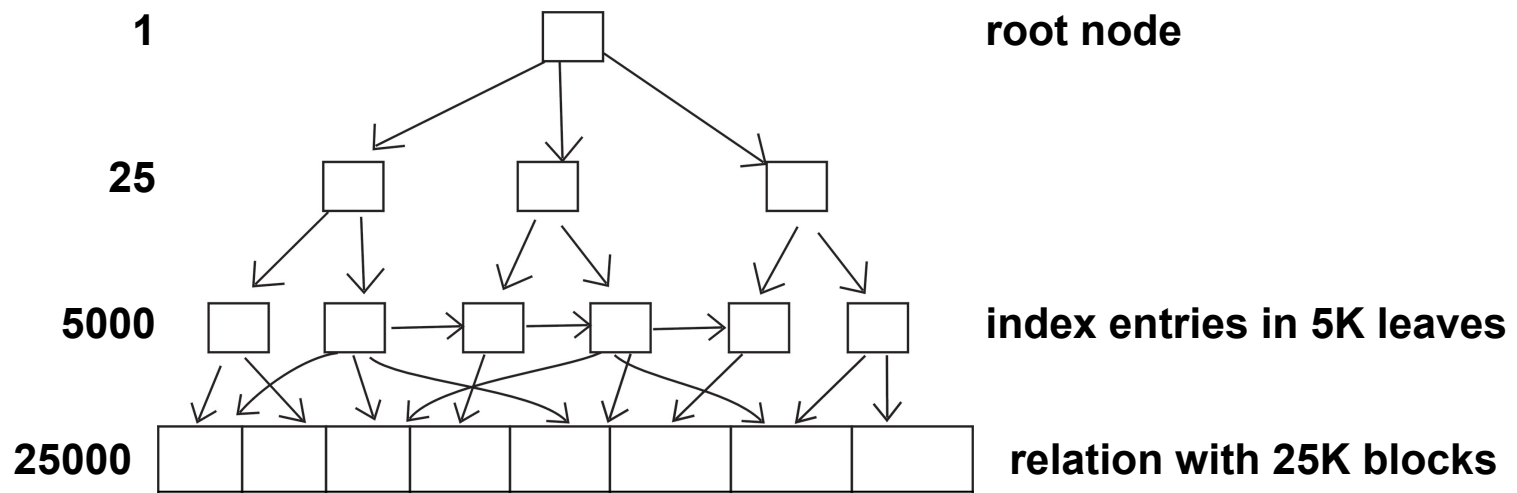


NYU

TANDON SCHOOL
OF ENGINEERING

Case 1: Unclustered B⁺-tree, 4KB blocks

- Key of size = 8 bytes and Pointer = 8 bytes
- → $n = 250$ (approximation)
- Assume average degree 200 (80%)



- Find one item: $4 * 10.4\text{ms} = 41.6\text{ms}$
- Range Query for 1000 records: $1000 * 10.4\text{ms} + 5 * 10.4\text{ms} + 2 * 10.4\text{ms} \approx 10.5\text{s}$



NEW YORK UNIVERSITY



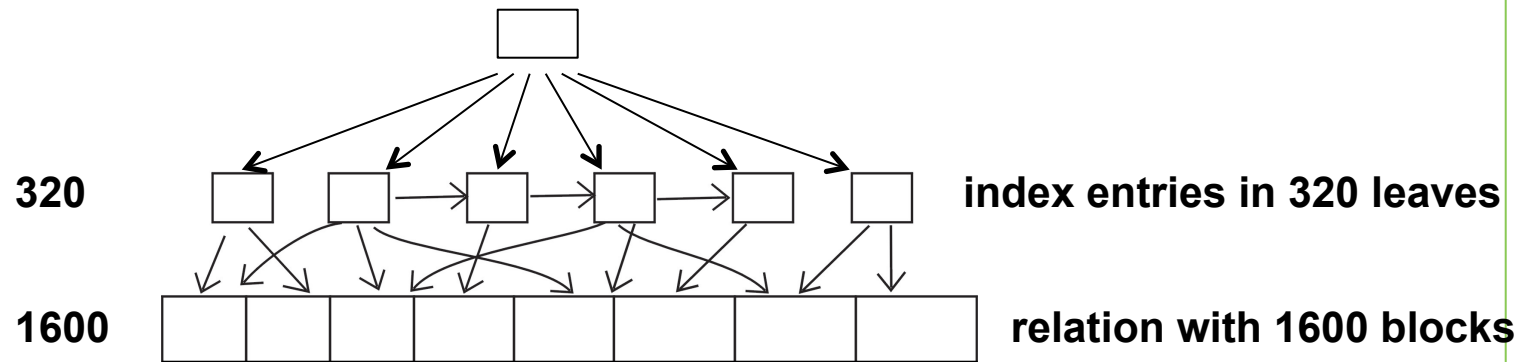
NYU

TANDON SCHOOL
OF ENGINEERING

Case 2: Unclustered with block size 64 KB

➔ ~ 1600 blocks with 640 records each

- $n = 4000$! and $T_B = 16.4\text{ms}$
- Assume average degree 3200



- Retrieve one item: $3 * 16.4\text{ms} = 49.2\text{ms}$
- Retrieve 1000 items: $1000 * 16.4\text{ms} + 2 * 16.4\text{ms} = 16.5\text{s}$
- Overall, large blocks not a good idea



NEW YORK UNIVERSITY

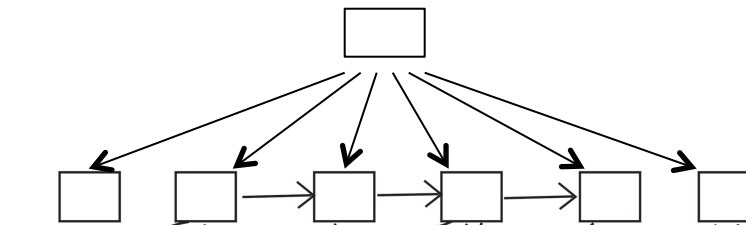


NYU

TANDON SCHOOL
OF ENGINEERING

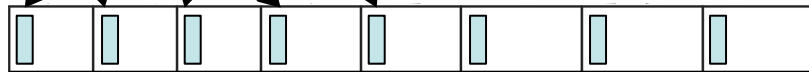
Case 3: Sparse clustered B+-Tree, 4KB blocks

- 125 blocks



index entries in leaves

- 25000 blocks



relation with blocks

- Retrieve one item $\rightarrow 3 * 10.4\text{ms} = 31.2\text{ms}$
- 1000 consecutive elements (range search) $\rightarrow (25+2) * 10.4 = 280\text{ms}$



NEW YORK UNIVERSITY



Case 4: Scanning the entire table

- 4KB page, 25000 pages of data
- $25000 * 10.4\text{ms} \approx 280\text{s}$ (4KB block model)
- -----
- 64KB page, 1600 page of data
- $1600 * 16.4\text{ms} \approx 26\text{s}$ (64KB block model)
- -----
- Scan model
- $10\text{ms} + 10 \text{ sec} = 10.01 \text{ sec}$ (To scan 100MB file)

➔ Disk model predicts that sometimes, scanning is better than using an unclustered index. Also, 64KB blocks worse





■ Composite Keys

- σ salary = 10000 and age = 30 (employees)
- σ salary > 80000 and 20 < age < 30 (employees)
- How to access these tuples:
 - Use index on salary
 - Use index on age
 - Use composite index
 - - (age, salary) or (salary, age)
 - - hash age \times salary
 - - 2 – D data structures (grid file)





■ Multi-Attribute B+-Trees

- **(age, salary) → 6 bytes per key**
 - age size is 2 bytes and salary size is 4 bytes
- **“Sort by age, and by salary if same age”**
 - **Age = 20 \wedge salary = 10000 ✓**
 - **Age = 20 \wedge salary > 10000 ✓**
 - **Age > 20 \wedge salary = 10000 ?**
 - **Age > 20 \wedge salary > 10000 ?**





NYU

TANDON SCHOOL
OF ENGINEERING

age	salary
20	18000
22	8000
22	11000
22	25000
22	40000
22	80000
24	11000
24	25000
25	25000
26	11000

salary	age
8000	22
11000	22
11000	24
11000	26
18000	20
25000	22
25000	24
25000	25
40000	22
80000	22



NEW YORK UNIVERSITY

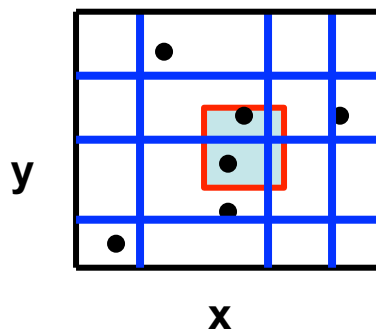
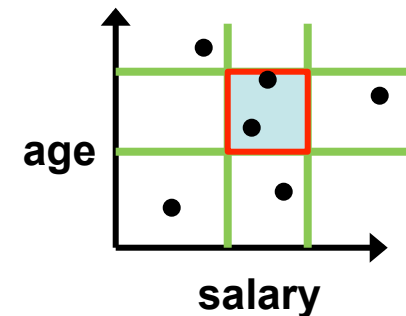


NYU

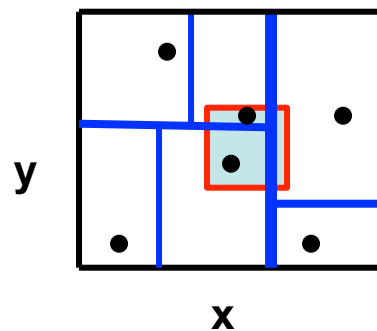
TANDON SCHOOL
OF ENGINEERING

Multi-Dimensional Indexes

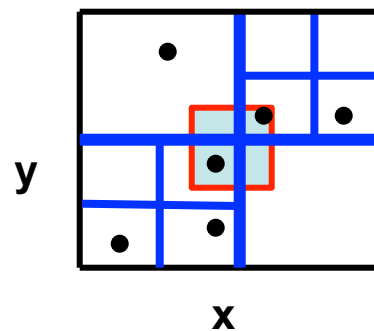
- Not same as multi-attribute composite
- Assume 2 attributes (dimensions)
- Often used in spatial DB and GIS
- Some examples:



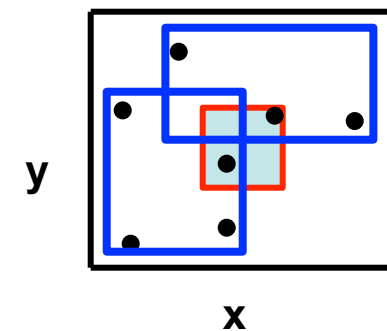
(a) grid file



(b) k-d tree (d=2)



(c) quad-tree



(d) R*-tree



NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

■ Bulk-Building B⁺-Trees

- Determine desired occupancy (90%)
- Extract (key, pointer) pairs from relation (scan)
- Sort pairs



- Build leaf nodes by scanning pairs in sorted order until node 90% full
- Repeat on next level
- Much much faster than inserting one item at a time



NEW YORK UNIVERSITY



NYU

**TANDON SCHOOL
OF ENGINEERING**



--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--

database table with 10M rows

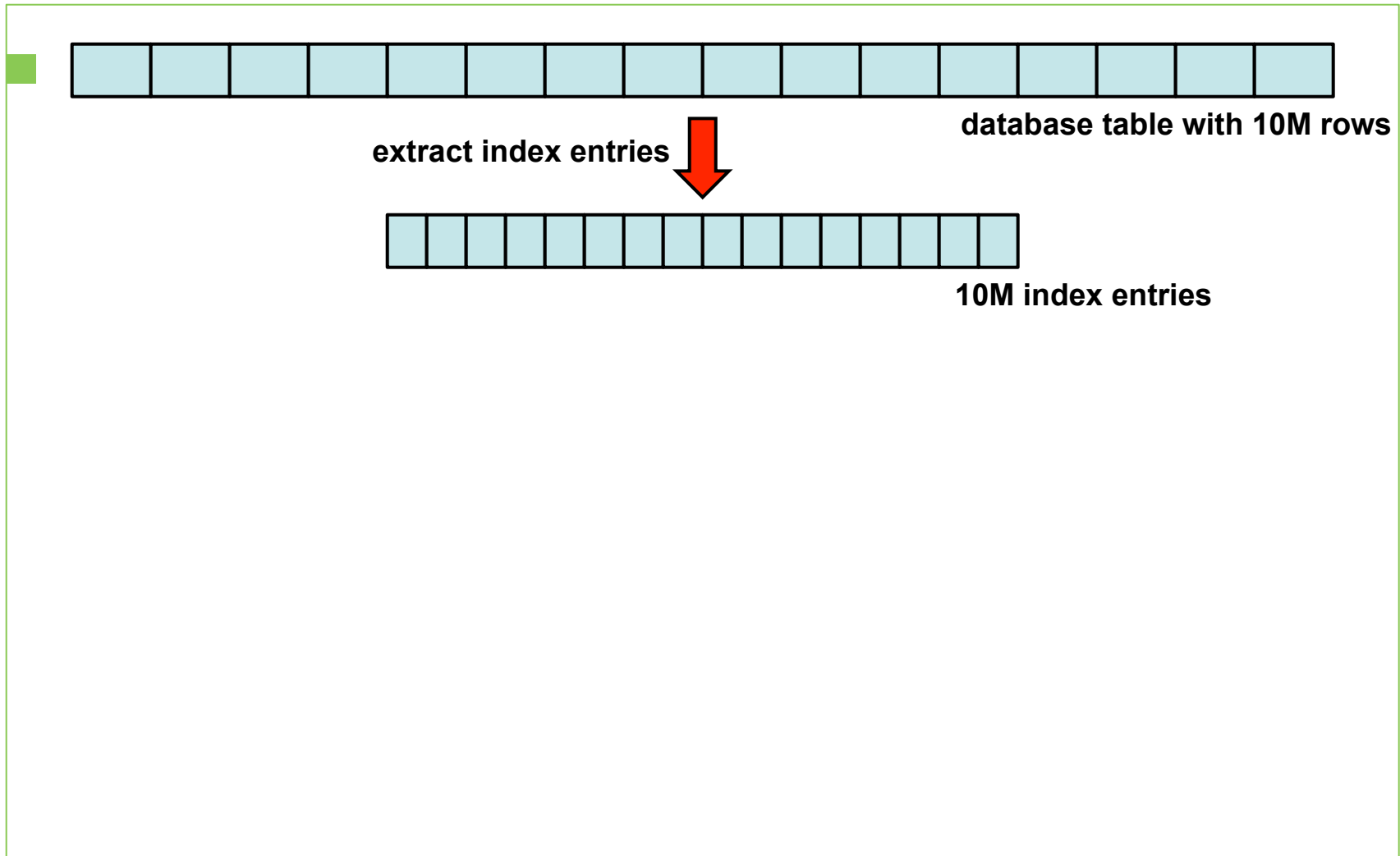


NEW YORK UNIVERSITY



NYU

**TANDON SCHOOL
OF ENGINEERING**

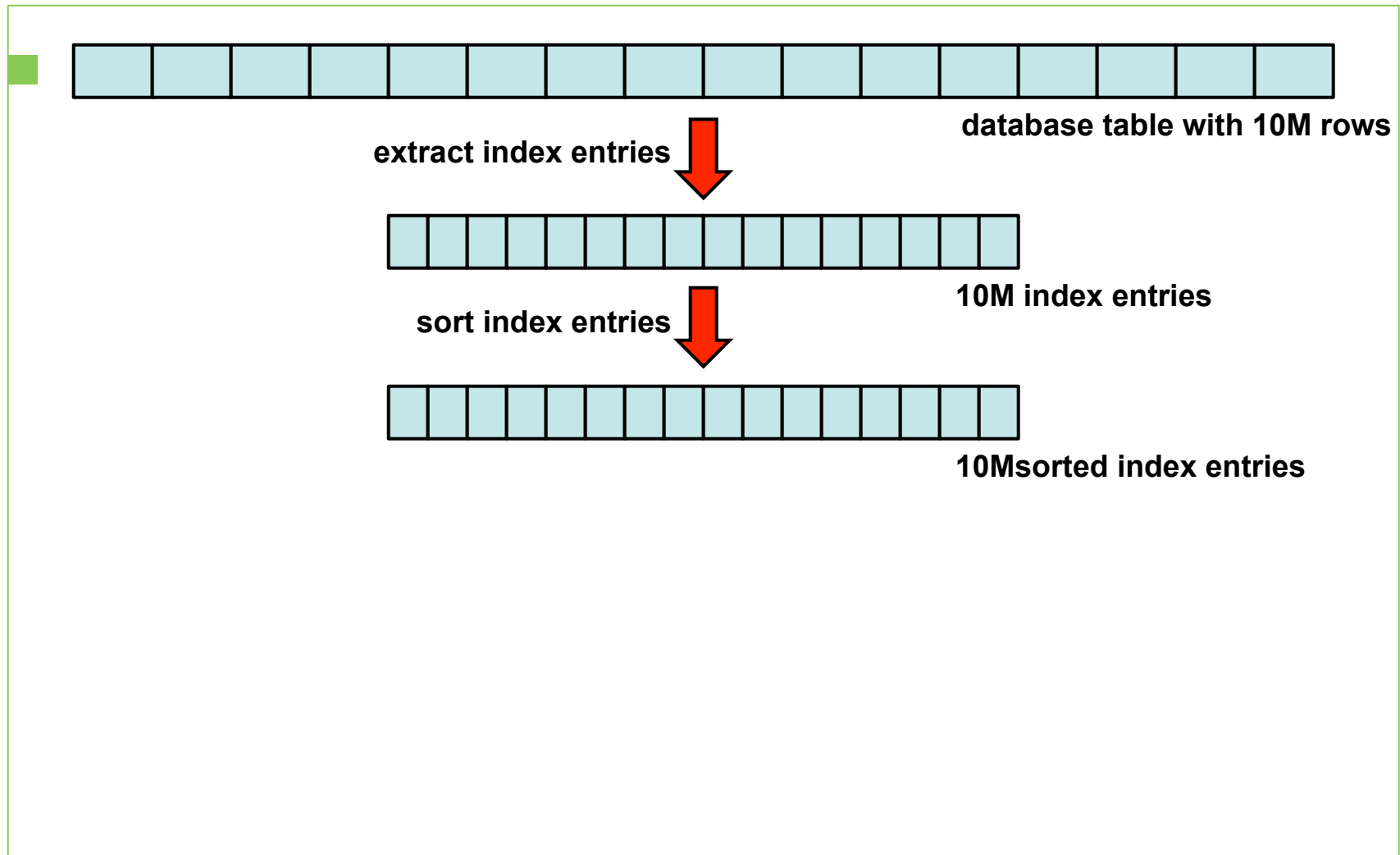


NEW YORK UNIVERSITY



NYU

**TANDON SCHOOL
OF ENGINEERING**

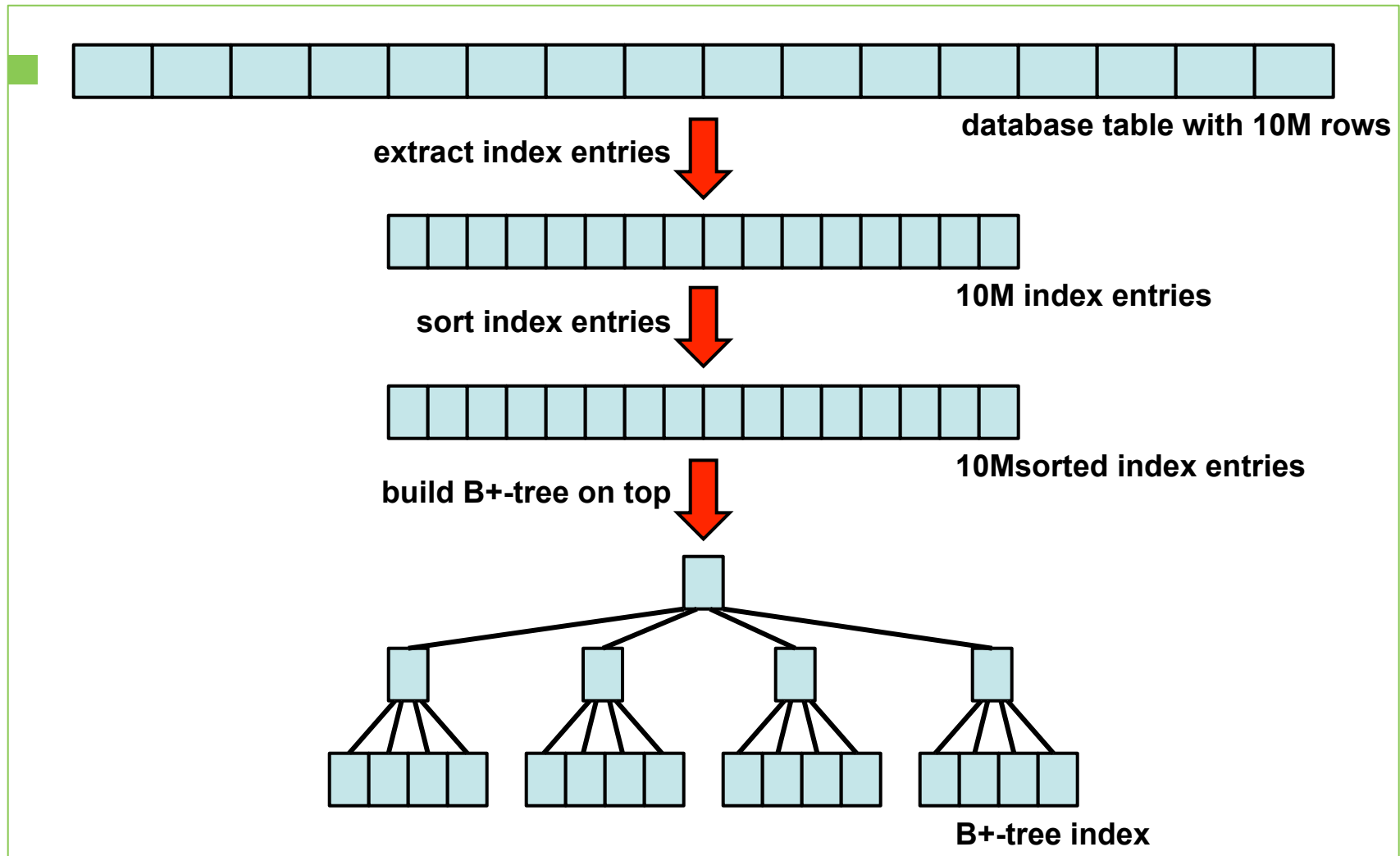


NEW YORK UNIVERSITY



NYU

**TANDON SCHOOL
OF ENGINEERING**

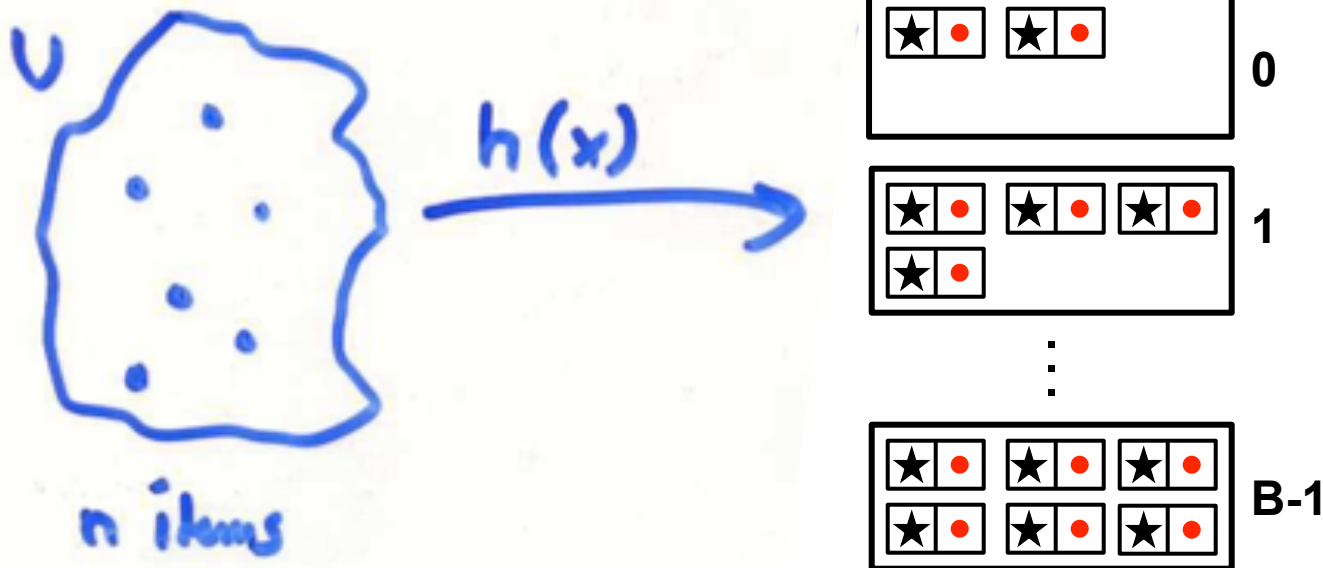


NEW YORK UNIVERSITY



■ Hashing Basics

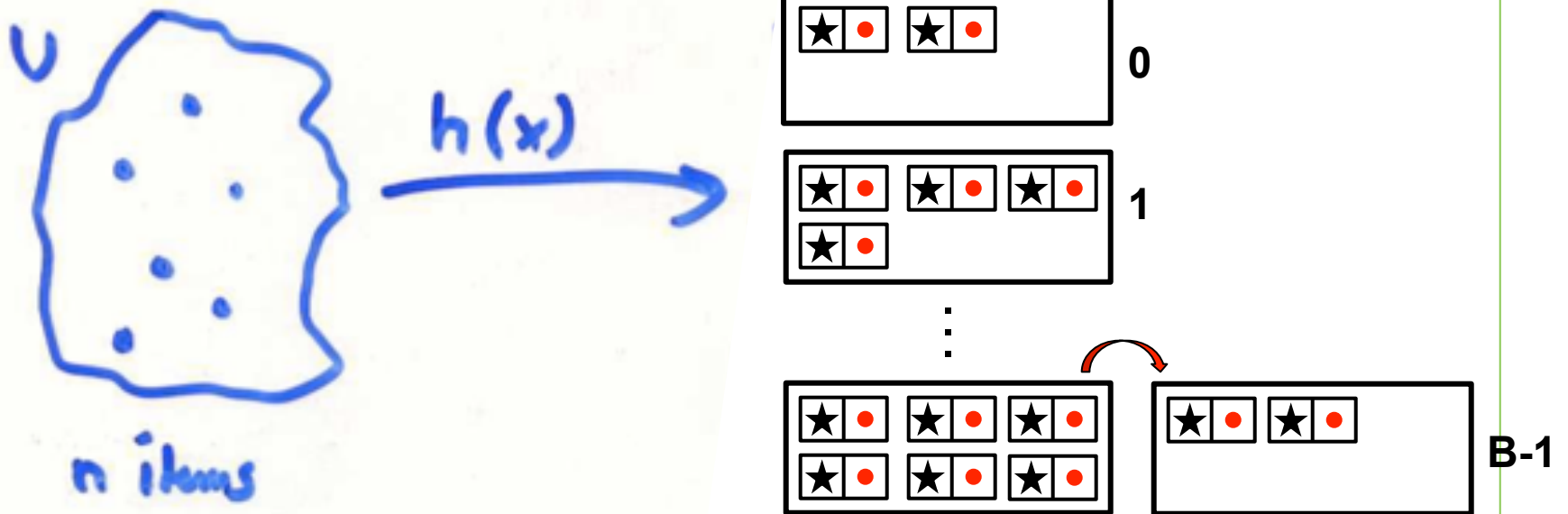
- Use hash function to assign index entries in to buckets
- If bucket \approx block then great !
- K = key domain, B = # buckets, $h: K \rightarrow [0, 1, 2, \dots, B-1]$





■ Problem: Some buckets may overflow

- Why does this happen?
- Static hashing: must add overflow blocks





NYU

TANDON SCHOOL
OF ENGINEERING

■ Extendible Hashing

- Hash $h: K \rightarrow [0, 1, 2, \dots, 2^{32}-1]$ (or even more bits)
- But at first only look at a few bits

10100101101110011001 ...

01101010010110110101 ...

01011010010110100101 ...

10010100101101111011 ...

10110101101110101001 ...



NEW YORK UNIVERSITY



Extendible Hashing

- Hash $h: K \rightarrow [0, 1, 2, \dots, 2^{32}-1]$ (or even more bits)
- But at first only look at a few bits

1	0	1	0	0	1	0	1	1	0	1	1	1	0	0	1	1	0	0	1	...
0	1	1	0	1	0	1	0	0	1	0	1	1	0	1	1	0	1	0	1	...
0	1	0	1	1	0	1	0	0	1	0	1	1	0	1	0	0	1	0	1	...
1	0	0	1	0	1	0	0	1	0	1	1	0	1	1	1	0	1	1	...	
1	0	1	1	0	1	0	1	1	0	1	1	1	0	1	0	1	0	0	1	...

- Maybe only 16 buckets first \rightarrow 4 bits





Extendible Hashing

- Hash $h: K \rightarrow [0, 1, 2, \dots, 2^{32}-1]$ (or even more bits)
- But at first only look at a few bits

1	0	1	0	0	1	0	1	1	0	1	1	1	0	0	1	1	0	0	1	...
0	1	1	0	1	0	1	0	0	1	0	1	1	0	1	1	0	1	0	1	...
0	1	0	1	1	0	1	0	0	1	0	1	1	0	1	0	0	1	0	1	...
1	0	0	1	0	1	0	0	1	0	1	1	0	1	1	1	0	1	1	...	
1	0	1	1	0	1	0	1	1	0	1	1	1	0	1	0	1	0	0	1	...

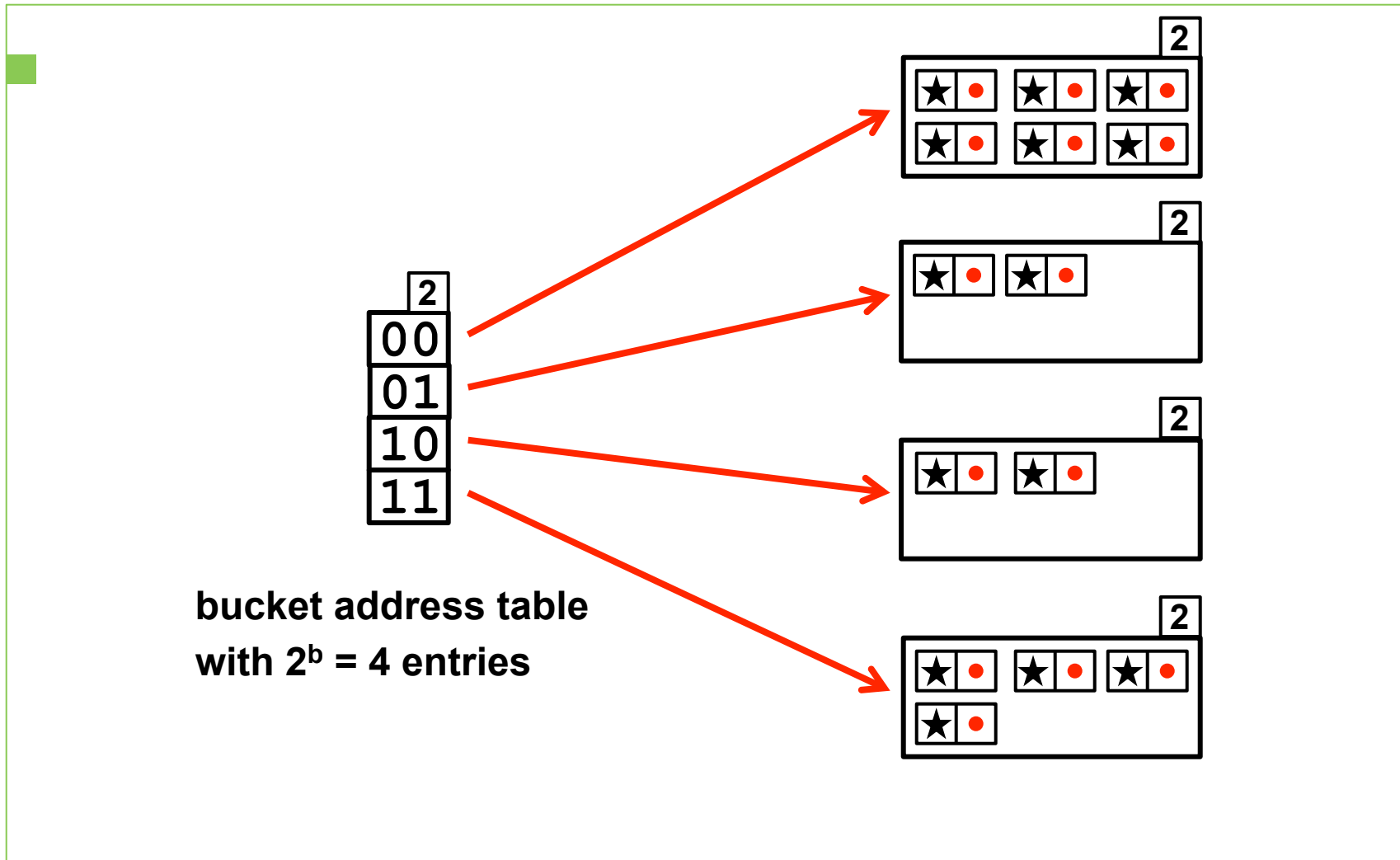
- Maybe only 16 buckets first \rightarrow 4 bits
- Later: look at more bits





NYU

TANDON SCHOOL
OF ENGINEERING

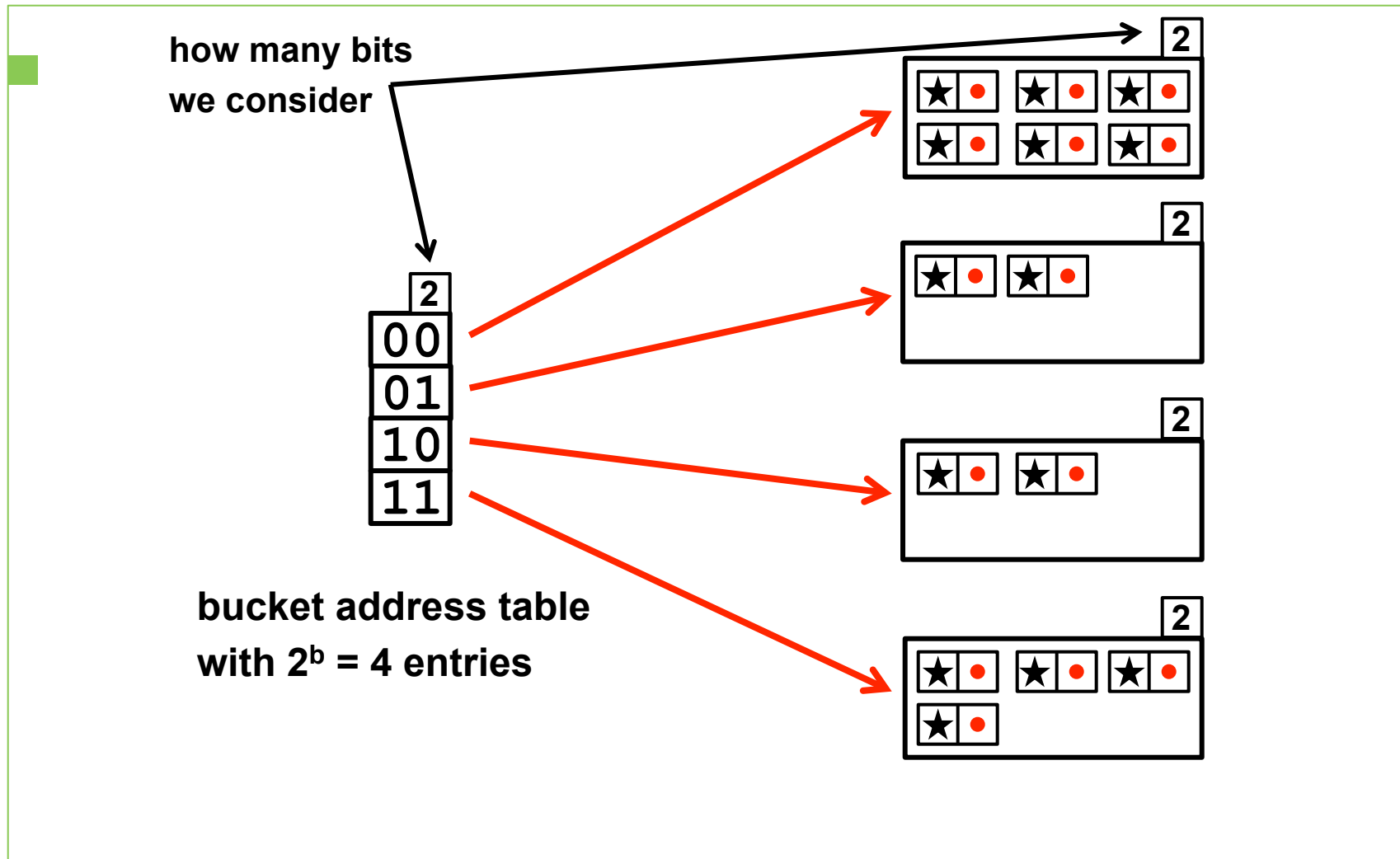


NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING



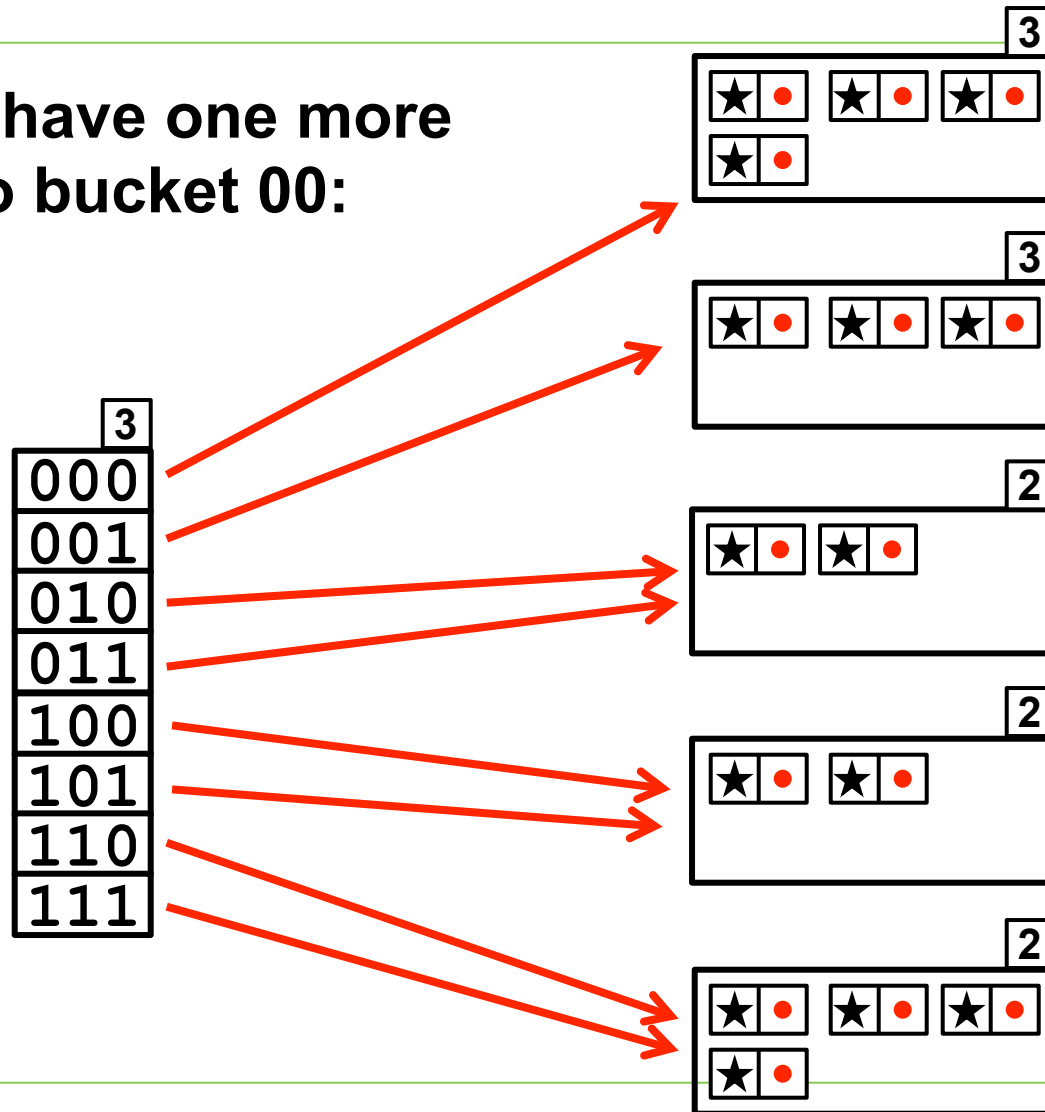
NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

Suppose we have one more
insertion into bucket 00:



NEW YORK UNIVERSITY



NYU

TANDON SCHOOL
OF ENGINEERING

■ **Extendible Hashing**

- **Uses bucket address table (BAT)**
- **No need to split non-full buckets**
- **BAT should fit in main memory**
- **Note: BAT about size of leaf level of dense B+-tree**

- **Cost to find one record: 2 seeks (bucket, record)**
- **Still need overflow buckets**



NEW YORK UNIVERSITY