# Order-Preserving Trees

Lecture #09

**Database Systems**
15-445/15-645
Fall 2017

**Andy Pavlo**
Computer Science Dept.
Carnegie Mellon Univ.

# ADMINISTRIVIA

**Project #1** is due Monday
October 2nd @ 11:59pm

**Homework #3** is due Wednesday
October 4th @ 11:59pm

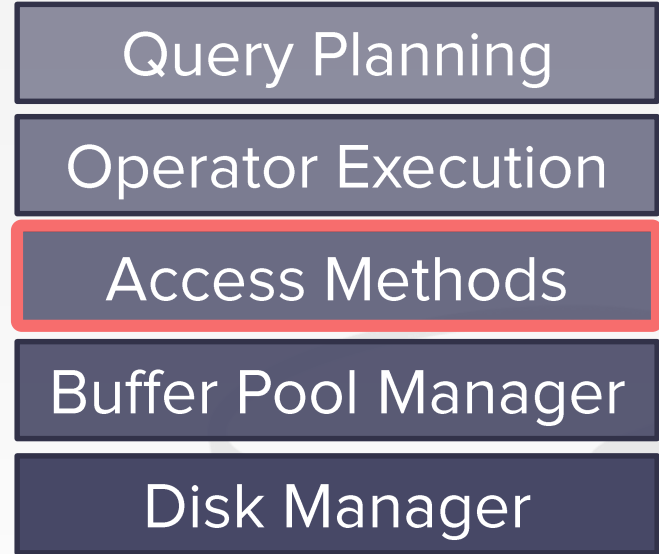CARNEGIE MELLON
**DATABASE GROUP**

# STATUS

We are now going to talk about how to support the DBMS's execution engine to read/write data from pages.

Two types of data structures:
→ Hash Tables
→ Trees

Query Planning

Operator Execution

Access Methods

Buffer Pool Manager

Disk Manager

CARNEGIE MELLON
DATABASE GROUP

# DATA STRUCTURES

Internal Meta-data

Core Data Storage

Temporary Data Structures

Table Indexes

CARNEGIE MELLON
DATABASE GROUP

# TABLE INDEXES

A **table index** is a replica of a subset of a table's columns that are organized and/or sorted for efficient access using a subset of those columns.

The DBMS ensures that the contents of the table and the index are always in sync.

# TABLE INDEXES

It is the DBMS's job to figure out the best index(es) to use to execute each query.

There is a trade-off on the number of indexes to create per database.
→ Storage Overhead
→ Maintenance Overhead

# TODAY'S AGENDA

B+Tree

Skip List

Radix Tree

Extra Index Stuff

CARNEGIE MELLON
DATABASE GROUP

# B-TREE FAMILY

There is a specific data structure called a **B-Tree**, but then people also use the term to generally refer to a class of data structures.

→ **B-Tree**
→ **B+Tree**
→ **B$^{link}$-Tree**
→ **B\*Tree**

CARNEGIE MELLON
**DATABASE GROUP**

# B+TREE

A **B+Tree** is a self-balancing tree data structure that keeps data sorted and allows searches, sequential access, insertions, and deletions in **O(log n)**.

→ Generalization of a binary search tree in that a node can have more than two children.

→ Optimized for systems that read and write large blocks of data.

**The Ubiquitous B-Tree**

DOUGLAS COMER

*Computer Science Department, Purdue University, West Lafayette, Indiana 47907*

B-trees have become, de facto, a standard for file organization. File indexes of users, dedicated database systems, and general-purpose access methods have all been proposed and implemented using B-trees. This paper reviews B-trees and shows why they have been so successful. It discusses the major variations of the B-tree, especially the B*-tree, contrasting the relative merits and costs of each implementation. It illustrates a general purpose access method which uses a B-tree.

*Keywords and Phrases:* B-tree, B*-tree, B*-tree, file organization, index

*CR Categories:* 3.73 3.74 4.33 4 34

**INTRODUCTION**

The secondary storage facilities available on large computer systems allow users to store, update, and recall data from large collections of information called files. A computer must retrieve an item and place it in main memory before it can be processed. In order to make good use of the computer resources, one must organize files intelligently, making the retrieval process efficient.

The choice of a good file organization depends on the kinds of retrieval to be performed. There are two broad classes of retrieval commands which can be illustrated by the following examples:

Sequential: "From our employee file, prepare a list of all employees' names and addresses," and

Random: "From our employee file, extract the information about employee J. Smith".

We can imagine a filing cabinet with three drawers of folders, one folder for each employee. The drawers might be labeled "A–G," "H–R," and "S–Z," while the folders

might be labeled with the employees' last names. A sequential request requires the searcher to examine the entire file, one folder at a time. On the other hand, a random request implies that the searcher, guided by the labels on the drawers and folders, need only extract one folder.

Associated with a large, randomly accessed file in a computer system is an *index* which, like the labels on the drawers and folders of the file cabinet, speeds retrieval by directing the searcher to the small part of the file containing the desired item. Figure 1 depicts a file and its index. An index may be physically integrated with the file, like the labels on employee folders, or physically separate, like the labels on the drawers. Usually the index itself is a file. If the index file is large, another index may be built on top of it to speed retrieval further, and so on. The resulting hierarchy is similar to the employee file, where the topmost index consists of labels on drawers, and the next level of index consists of labels on folders.

Natural hierarchies, like the one formed by considering last names as index entries, do not always produce the best perform-

Computing Surveys, Vol 11, No 2, June 1979

CARNEGIE MELLON
**DATABASE GROUP**

# B+TREE: PROPERTIES

A B+tree is an *M*-way search tree
with the following properties:
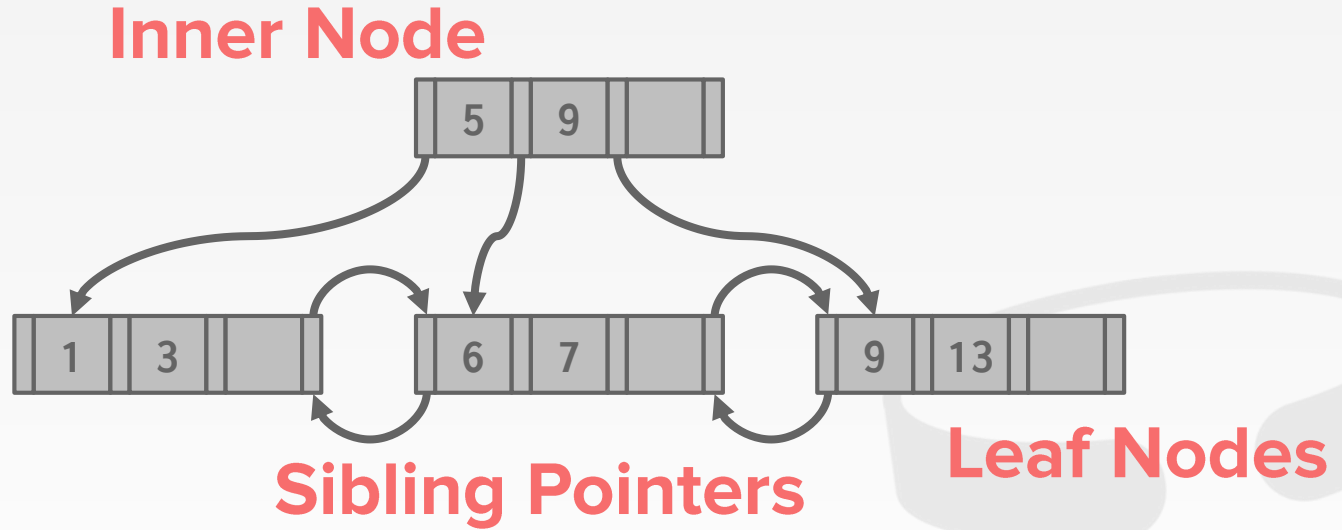→ It is perfectly balanced (i.e., every leaf
  node is at the same depth).
→ Every inner node other than the root,
  is at least half-full
  $M/2-1 \leq \#keys \leq M-1$
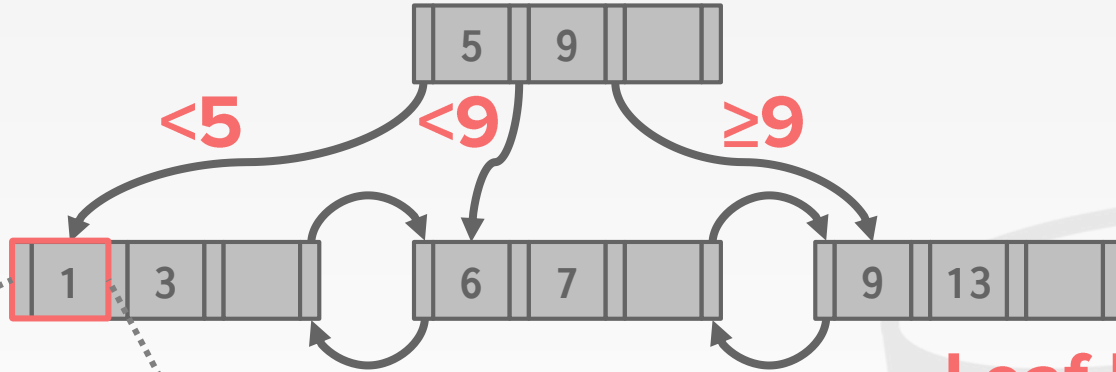→ Every inner node with **k** keys has **k+1**
  non-null children

# B+TREE OVERVIEW



Inner Node

Sibling Pointers

Leaf Nodes

CARNEGIE MELLON
DATABASE GROUP

# B+TREE OVERVIEW



Inner Node

5  9

<5    <9    ≥9

1  3    6  7    9  13

Sibling Pointers

Leaf Nodes

<value>/<key>

# B+TREE NODES

Every node in the B+Tree contains an array of key/value pairs.
→ The keys will always be the column or columns that you built your index on
→ The values will differ based on whether the node is classified as **inner nodes** or **leaf nodes.**

The arrays are always kept in sorted order.

CARNEGIE MELLON
**DATABASE GROUP**

# B+TREE: LEAF NODE VALUES

## Approach #1: Record Ids
→ A pointer to the location of the tuple that the index entry corresponds to.

## Approach #2: Tuple Data
→ The actual contents of the tuple is stored in the leaf node.
→ Secondary indexes have to store the record id as their values.

CARNEGIE MELLON
DATABASE GROUP

# B+TREE LEAF NODES



B+Tree Leaf Node

# B+TREE LEAF NODES



B+Tree Leaf Node

Prev
Next
PageID ← ¤ | K1 | ¤ • • • Kn | ¤ | ¤ → PageID

Key+Value

# B+TREE LEAF NODES



**B+Tree Leaf Node**

| Level | Slots | Prev | Next |
|-------|-------|------|------|
| # | # | ¤ | ¤ |

**Sorted Keys**

| K1 | K2 | K3 | K4 | K5 | • • • | Kn |

**Values**

| ¤ | ¤ | ¤ | ¤ | ¤ | • • • | ¤ |

# B-TREE VS. B+TREE

The original **B-Tree** from 1972 stored keys + values in all nodes in the tree.
→ More space efficient since each key only appears once in the tree.

A **B+Tree** only stores values in leaf nodes. Inner nodes only guide the search process.

# B+TREE: INSERT

Find correct leaf **L**.

Put data entry into **L** in sorted order.
→ If L has enough space, done!
→ Else, must split **L** into **L** and a new node **L2**
  • Redistribute entries evenly, copy up middle key.
  • Insert index entry pointing to **L2** into parent of **L**.

To split inner node, redistribute entries evenly, but push up middle key.

Source: Chris Re

CARNEGIE MELLON
DATABASE GROUP

# B+TREE VISUALIZATION

http://cmudb.io/btree

https://www.cs.usfca.edu/~gall es/visualization/BPlusTree.html

# B+TREE: DELETE

Start at root, find leaf **L** where entry belongs.

Remove the entry.
→ If **L** is at least half-full, done!
→ If **L** has only `M/2-1` entries,
  • Try to re-distribute, borrowing from sibling (adjacent node with same parent as **L**).
  • If re-distribution fails, merge **L** and sibling.

If merge occurred, must delete entry (pointing to **L** or sibling) from parent of **L**.

# B+TREES IN PRACTICE

Typical Fill-Factor: 67%.
→ Average Fanout = 2*100*0.67 = 134

Typical Capacities:
→ Height 4: 1334 = 312,900,721 entries
→ Height 3: 1333 =     2,406,104 entries

Pages per level:
→ Level 1  =          1 page   =     8 KB
→ Level 2 =      134 pages  =    1 MB
→ Level 3 =  17,956 pages = 140 MB

# B+TREE DESIGN CHOICES

Merge Threshold

Non-Unique Indexes

Variable Length Keys

Prefix Compression

CARNEGIE MELLON
DATABASE GROUP

# B+TREE: MERGE THRESHOLD

Some DBMSs don't always merge nodes when it is half full.

Delaying a merge operation may reduce the amount of reorganization.

CARNEGIE MELLON
**DATABASE GROUP**

# B+TREE: NON-UNIQUE INDEXES

**Approach #1: Duplicate Keys**
→ Use the same leaf node layout but store duplicate keys multiple times.

**Approach #2: Value Lists**
→ Store each key only once and maintain a linked list of unique values.

# B+TREE: DUPLICATE KEYS

## B+Tree Leaf Node

| Level | Slots | Prev | Next |
|-------|-------|------|------|
| # | # | ¤ | ¤ |

**Sorted Keys**

| K1 | K1 | K1 | K2 | K2 | • • • | Kn |
|----|----|----|----|----|-------|----|

**Values**

| ¤ | ¤ | ¤ | ¤ | ¤ | • • • | ¤ |
|---|---|---|---|---|-------|---|

CARNEGIE MELLON
**DATABASE GROUP**

# B+TREE: VALUE LISTS



**B+Tree Leaf Node**

# B+TREE: VARIABLE LENGTH KEYS

## Approach #1: Pointers
→ Store the keys as pointers to the tuple's attribute.

## Approach #2: Variable Length Nodes
→ The size of each node in the B+Tree can vary.
→ Requires careful memory management.

## Approach #3: Key Map
→ Embed an array of pointers that map to the key + value list within the node.

# B+TREE: PREFIX COMPRESSION

The keys in the inner nodes are only used to "direct traffic".
→ We don't actually need the entire key.

Store a minimum prefix that is needed to correctly route probes into the index.

# B+TREE: PREFIX COMPRESSION

The keys in the inner nodes are only used to "direct traffic".
→ We don't actually need the entire key.

Store a minimum prefix that is needed to correctly route probes into the index.

# B+TREE BULK INSERT

The fastest/best way to build a
B+Tree is to first sort the keys and
then build the index from the
bottom up.

CARNEGIE MELLON
DATABASE GROUP

# B+TREE BULK INSERT

The fastest/best way to build a B+Tree is to first sort the keys and then build the index from the bottom up.

*Keys: 3, 7, 9, 13, 6, 1*

*Sorted Keys: 1, 3, 6, 7, 9, 13*

# B+TREE BULK INSERT

The fastest/best way to build a B+Tree is to first sort the keys and then build the index from the bottom up.

*Keys: 3, 7, 9, 13, 6, 1*
*Sorted Keys: 1, 3, 6, 7, 9, 13*

# B+TREE BULK INSERT

The fastest/best way to build a B+Tree is to first sort the keys and then build the index from the bottom up.

*Keys: 3, 7, 9, 13, 6, 1*

*Sorted Keys: 1, 3, 6, 7, 9, 13*

# OBSERVATION

The easiest way to implement a **dynamic** order-preserving index is to use a sorted linked list.

All operations have to linear search.
→ Average Cost: O(N)

# OBSERVATION

The easiest way to implement a **dynamic** order-preserving index is to use a sorted linked list.

All operations have to linear search.
→ Average Cost: O(N)

CARNEGIE MELLON
**DATABASE GROUP**

# SKIP LISTS

Invented in 1990.

Multiple levels of linked lists with extra pointers that **skip** over intermediate nodes.

Maintains keys in sorted order without requiring global rebalancing.



Skip Lists: A Probabilistic Alternative to Balanced Trees

# SKIP LISTS

A collection of lists at different levels
→ Lowest level is a sorted, singly linked list of all keys
→ 2nd level links every other key
→ 3rd level links every fourth key
→ In general, a level has half the keys of one below it

To insert a new key, flip a coin to decide how many levels to add the new key into.
Provides approximate **O(log n)** search times.

# SKIP LISTS: EXAMPLE

# SKIP LISTS: EXAMPLE

# SKIP LISTS: EXAMPLE

# SKIP LISTS: EXAMPLE

**Levels**

**End**



P=N/4

P=N/2

K2

K4

∞

∞

P=N

K1
V1

K2
V2

K3
V3

K4
V4

K6
V6

∞

CARNEGIE MELLON
**DATABASE GROUP**

# SKIP LISTS: INSERT

## Insert K5

# SKIP LISTS: INSERT

## Insert K5

# SKIP LISTS: INSERT

*Insert K5*

# SKIP LISTS: INSERT

## Insert K5

# SKIP LISTS: SEARCH

## *Find K3*

# SKIP LISTS: SEARCH

### Find K3

**Levels**

**End**

$K3<K5$



P=N/4

P=N/2

P=N

K2

K2

K1
V1

K2
V2

K3
V3

K4

K4
V4

K5

K5

K5
V5

K6
V6

∞

∞

∞

CARNEGIE MELLON
DATABASE GROUP

# SKIP LISTS: SEARCH

*Find K3*

# SKIP LISTS: SEARCH

## Find K3

# SKIP LISTS: SEARCH

*Find K3*

# SKIP LISTS: DELETE

First **logically** remove a key from the index by setting a flag to tell threads to ignore.

Then **physically** remove the key once we know that no other thread is holding the reference.
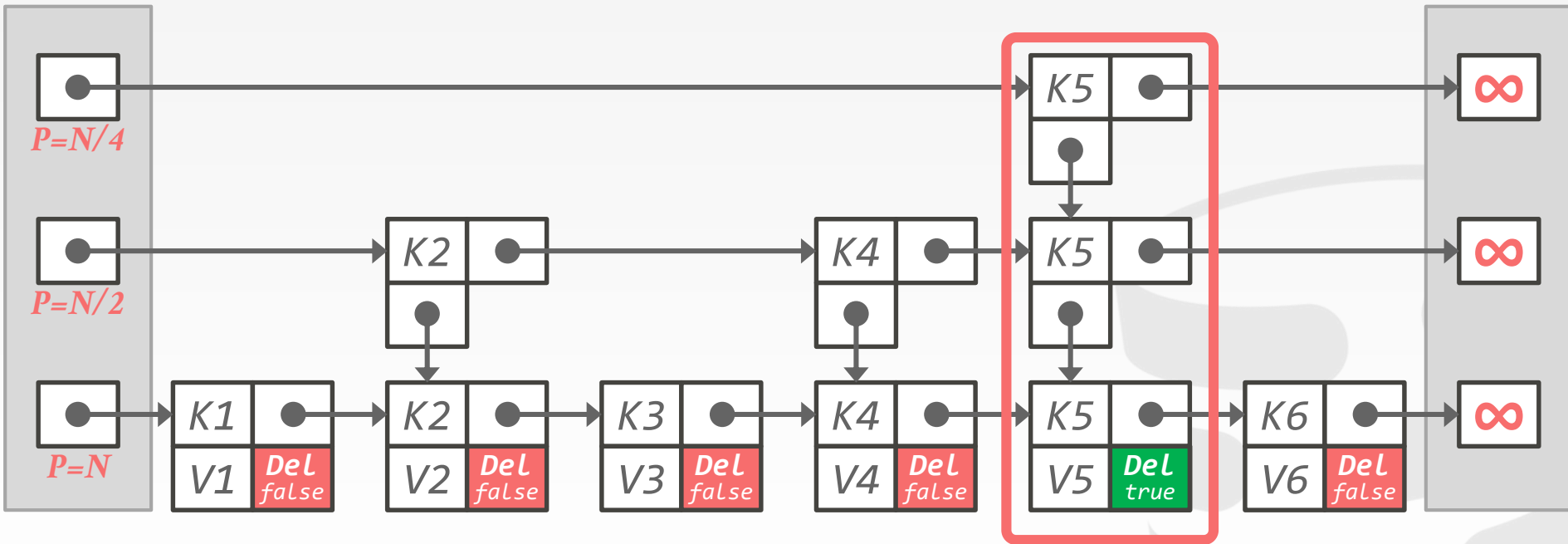
# SKIP LISTS: DELETE
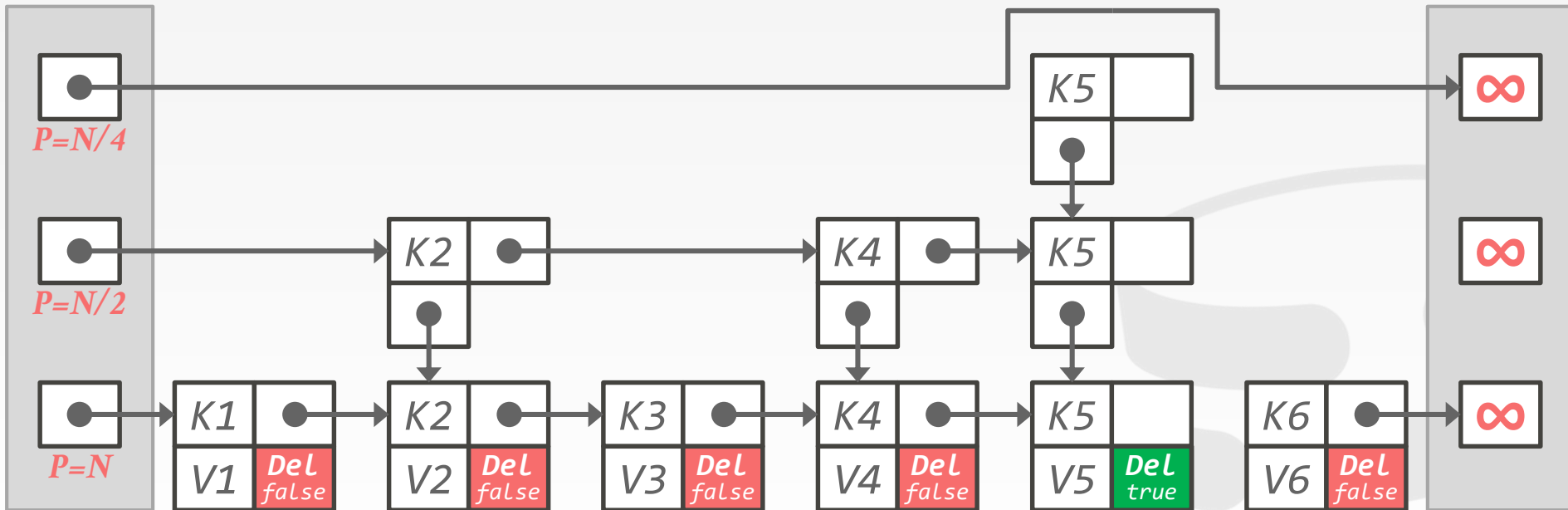
## Delete K5

# SKIP LISTS: DELETE



Delete K5

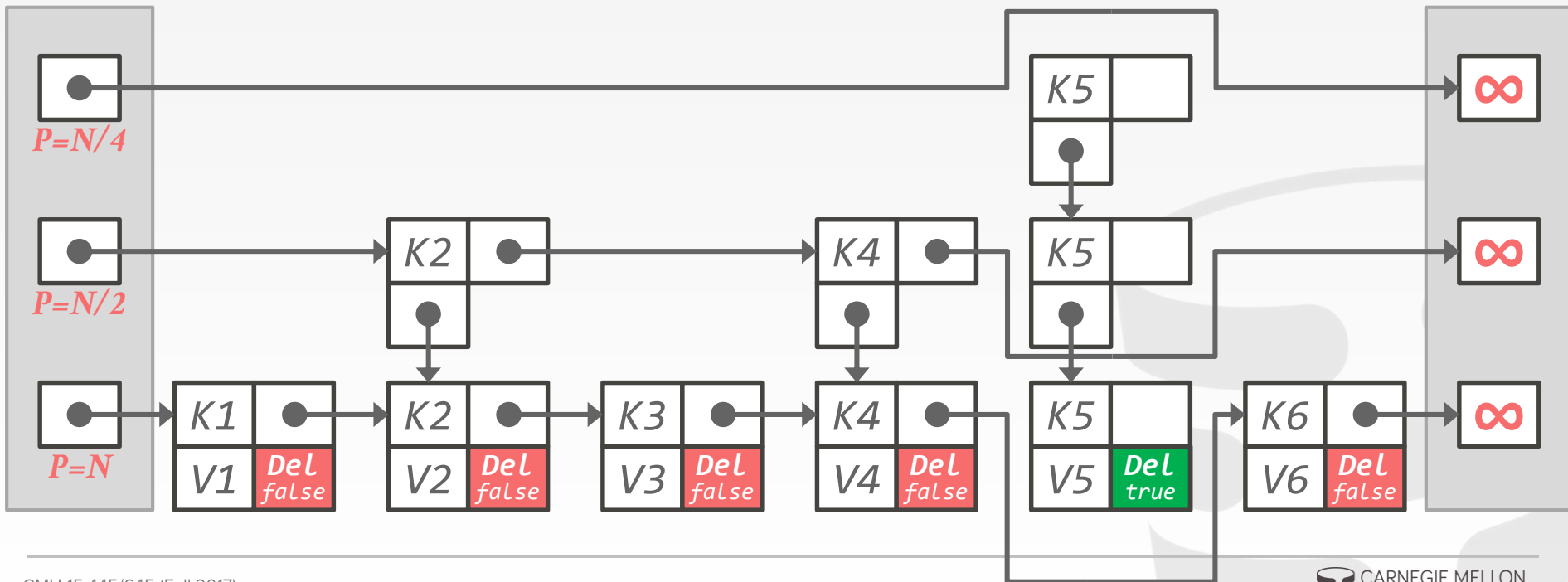# SKIP LISTS: DELETE

## Delete K5

# SKIP LISTS: DELETE
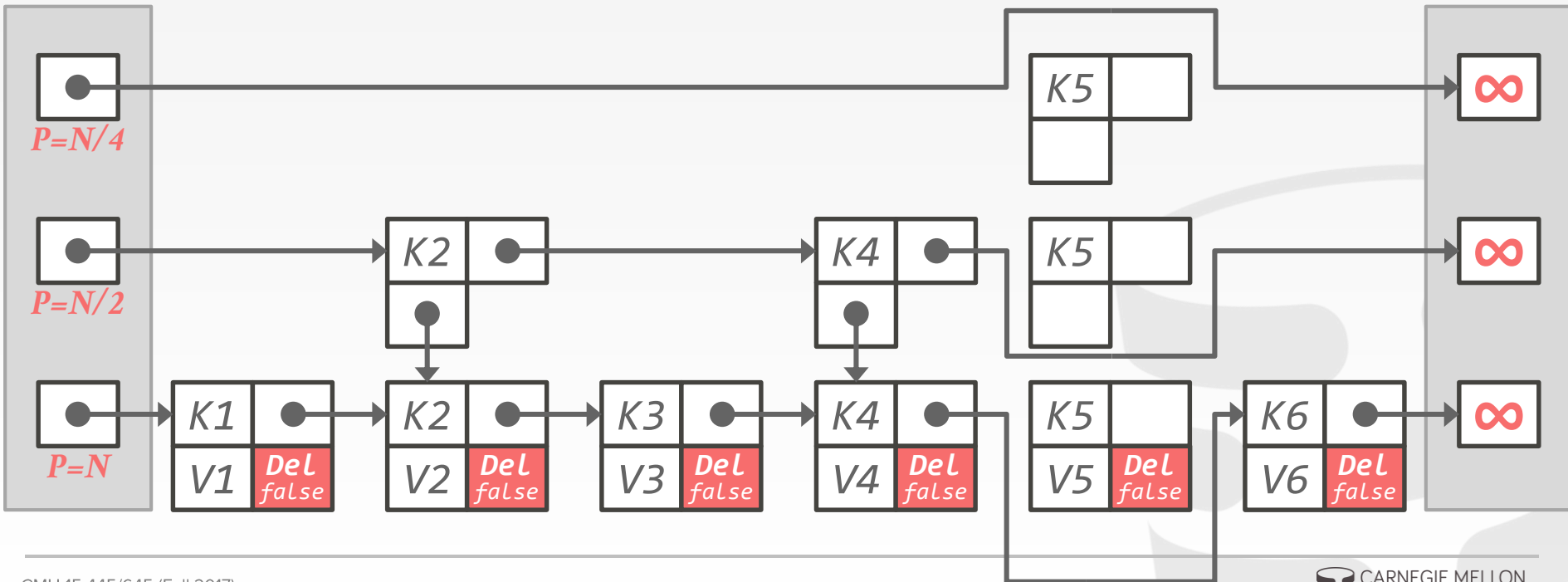


Delete K5

# SKIP LISTS: DELETE



Delete K5

# SKIP LISTS: DELETE

## Delete K5

# SKIP LISTS: ADVANTAGES

Uses less memory than a typical B+Tree
if you don't include reverse pointers.

Insertions and deletions do not require
rebalancing.

CARNEGIE MELLON
DATABASE GROUP

# SKIP LISTS: DISADVANTAGES

Not disk/cache friendly because they do not optimize locality of references.

Invoking random number generator multiple times per insert is slow.

Reverse search is non-trivial.

CARNEGIE MELLON
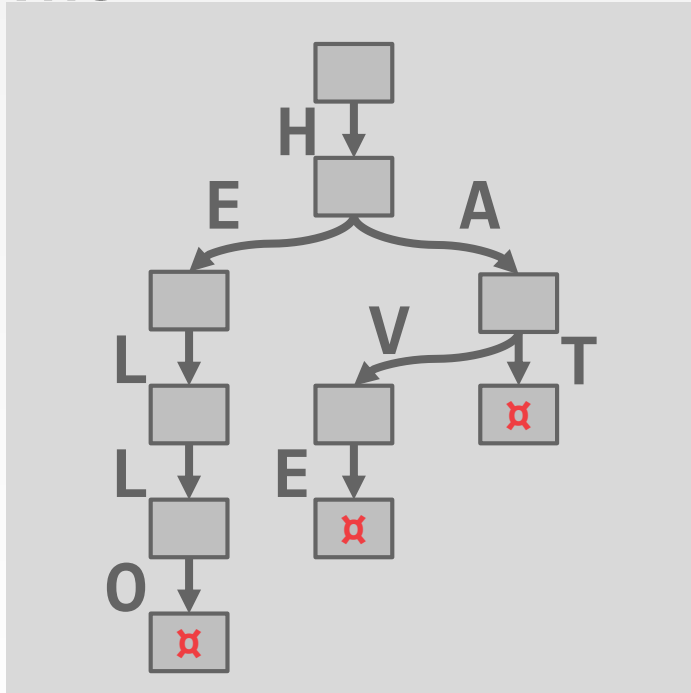DATABASE GROUP

# RADIX TREE

Uses digital representation of keys to examine prefixes one-by-one instead of comparing entire key.
→ The height of the tree depends on the length of keys.
→ Does not require rebalancing
→ The path to a leaf node represents the key of the leaf
→ Keys are stored implicitly and can be reconstructed from paths.
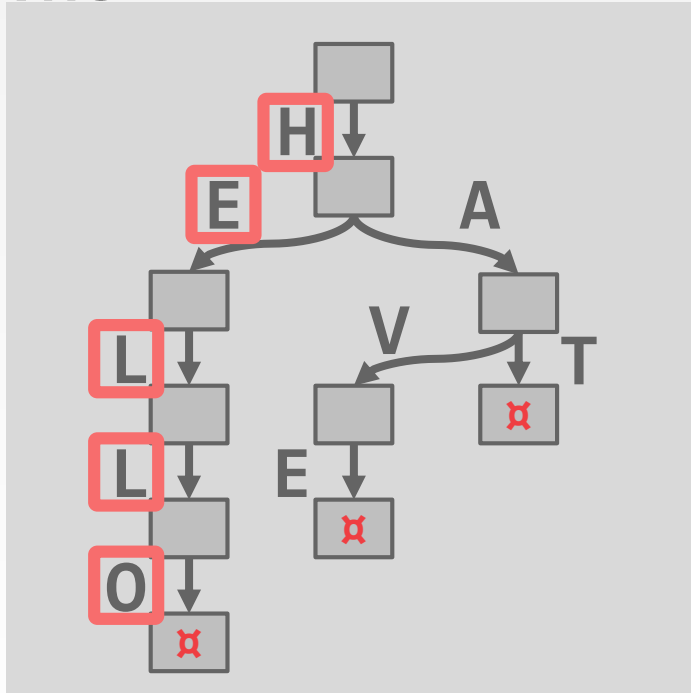
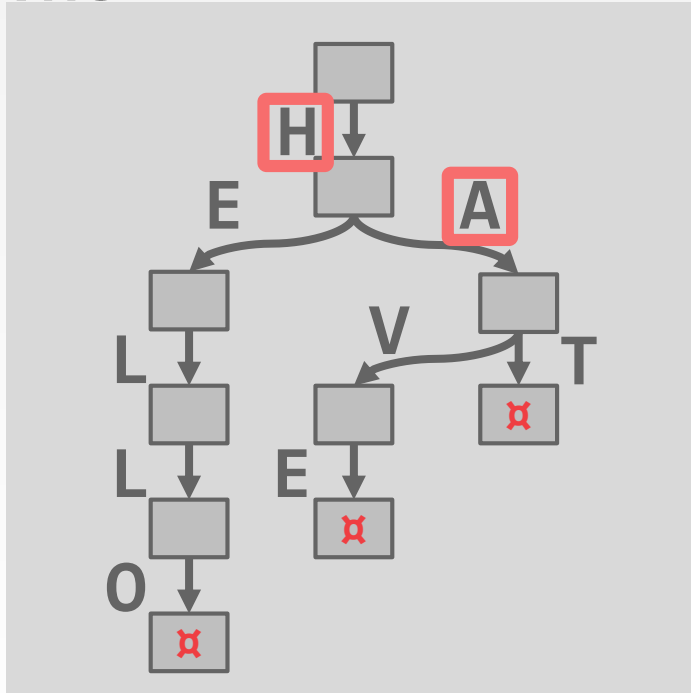# TRIE VS. RADIX TREE

*Trie*



Keys: HELLO, HAT, HAVE

CARNEGIE MELLON
DATABASE GROUP

# TRIE VS. RADIX TREE

*Trie*



Keys: HELLO HAT, HAVE

CARNEGIE MELLON
DATABASE GROUP

# TRIE VS. RADIX TREE
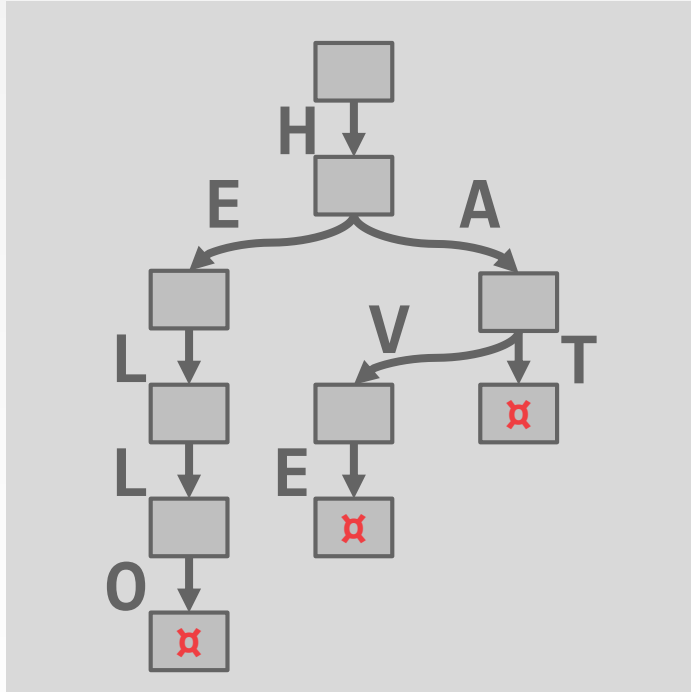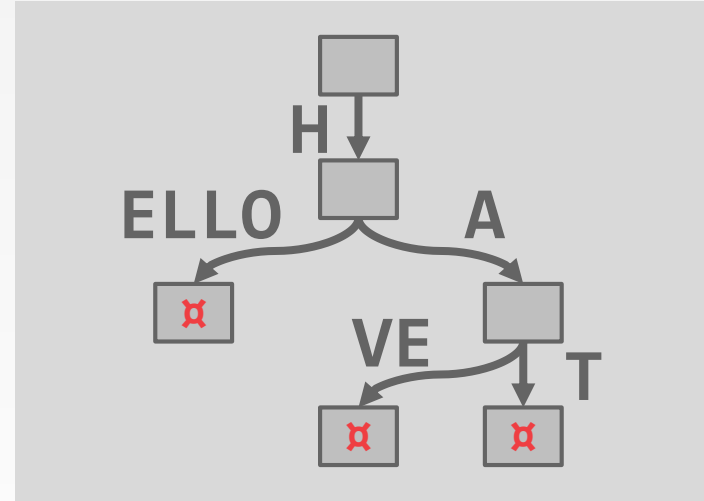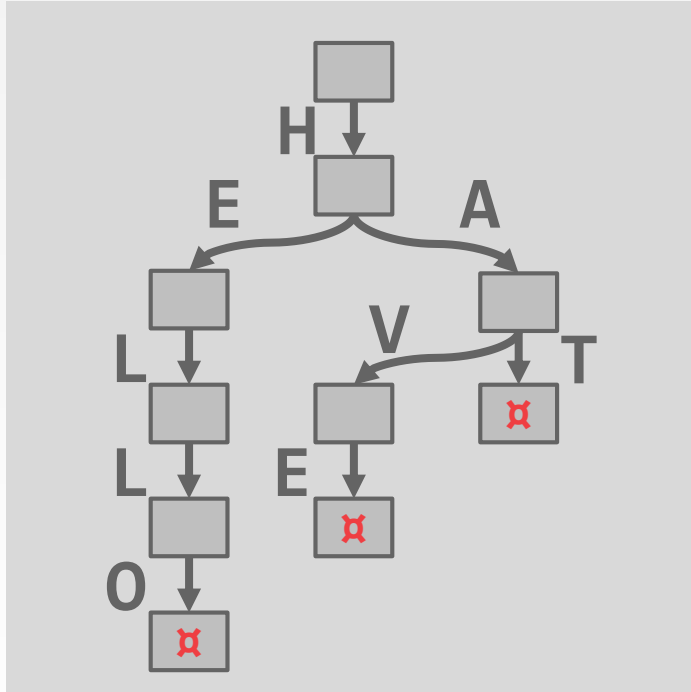
*Trie*



Keys: HELLO, HAT, HAVE

# TRIE VS. RADIX TREE



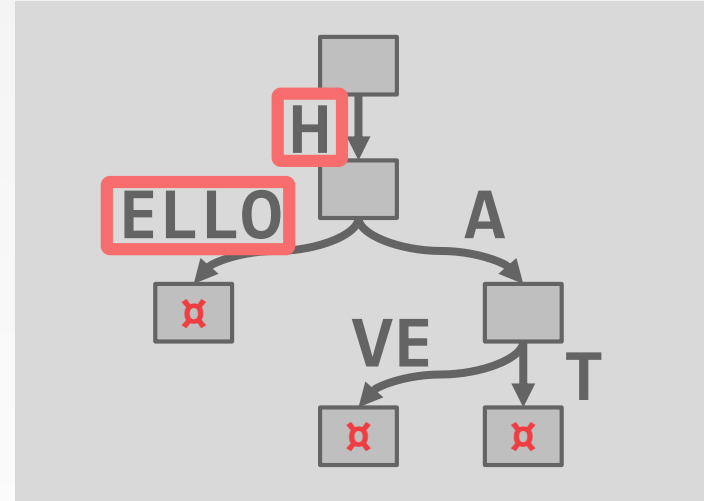Keys: HELLO, HAT, HAVE

# TRIE VS. RADIX TREE
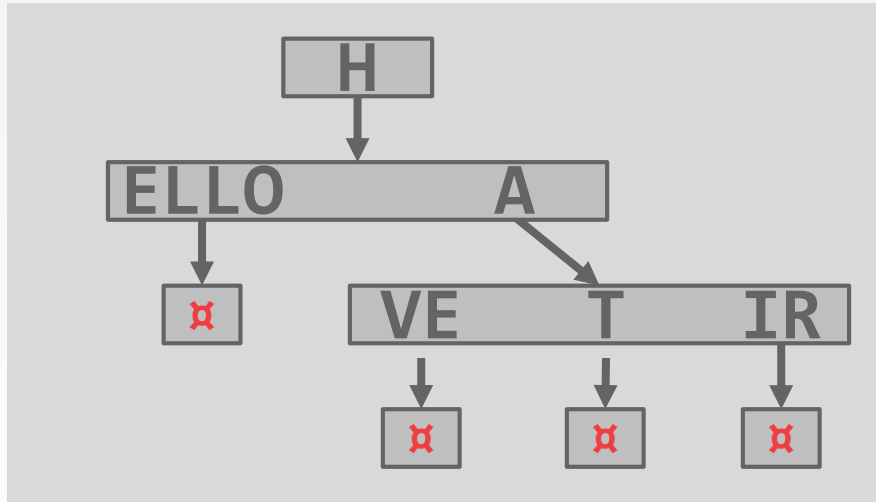


*Trie*

*Radix Tree*

Keys: HELLO HAT, HAVE

CARNEGIE MELLON
DATABASE GROUP
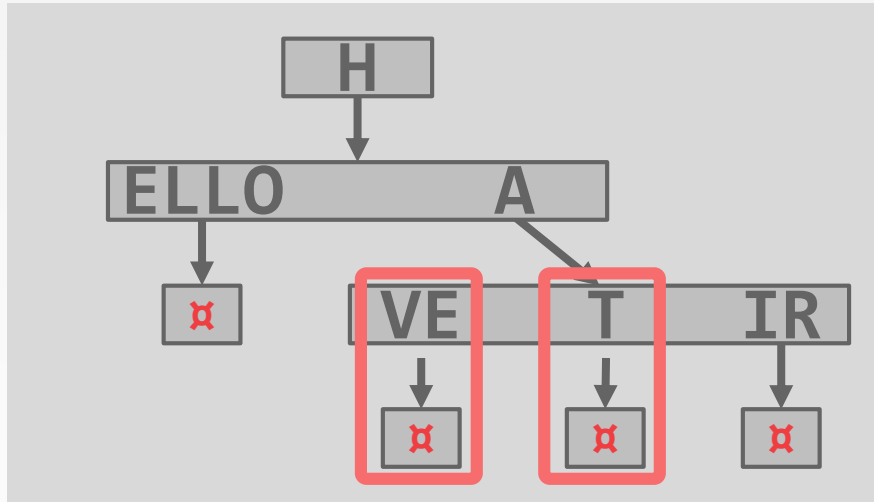
# RADIX TREE: MODIFICATIONS

# RADIX TREE: MODIFICATIONS



**Insert HAIR**

# RADIX TREE: MODIFICATIONS



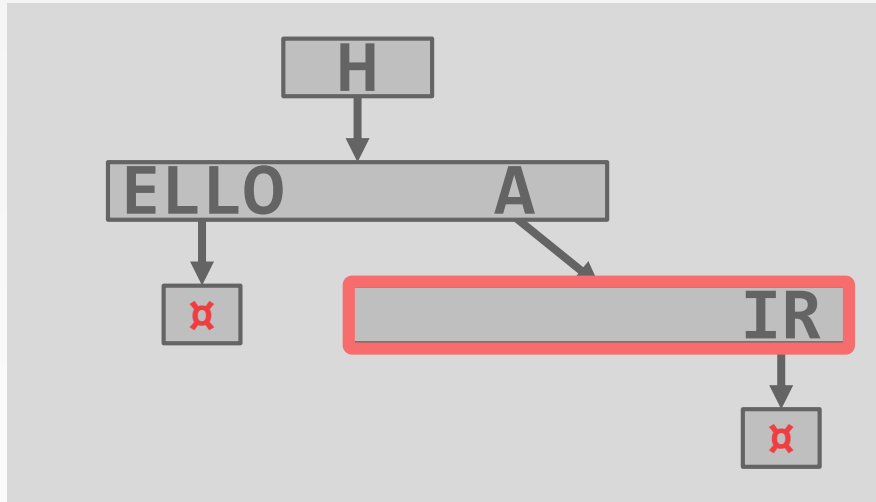**Insert HAIR**

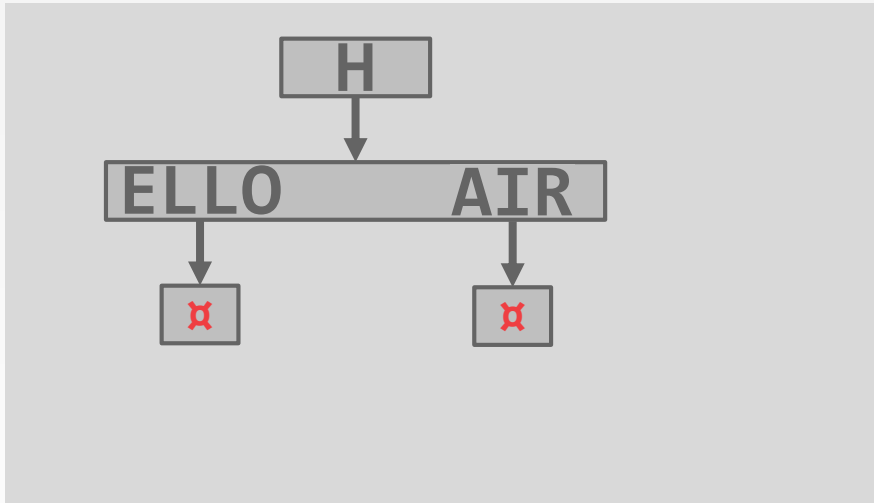**Delete HAT, HAVE**

# RADIX TREE: MODIFICATIONS



**Insert HAIR**

**Delete HAT, HAVE**

Insert HAIR

Delete HAT, HAVE

# RADIX TREE: MODIFICATIONS



**Insert HAIR**

**Delete HAT, HAVE**

CARNEGIE MELLON
**DATABASE GROUP**

# RADIX TREE: BINARY COMPARABLE KEYS

Not all attribute types can be decomposed into binary comparable digits for a radix tree.
→ **Unsigned Integers:** Byte order must be flipped for little endian machines.
→ **Signed Integers:** Flip two's-complement so that negative numbers are smaller than positive.
→ **Floats:** Classify into group (neg vs. pos, normalized vs. denormalized), then store as unsigned integer.
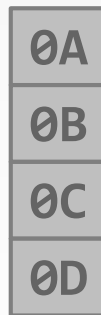→ **Compound:** Transform each attribute separately.

# RADIX TREE: BINARY COMPARABLE KEYS

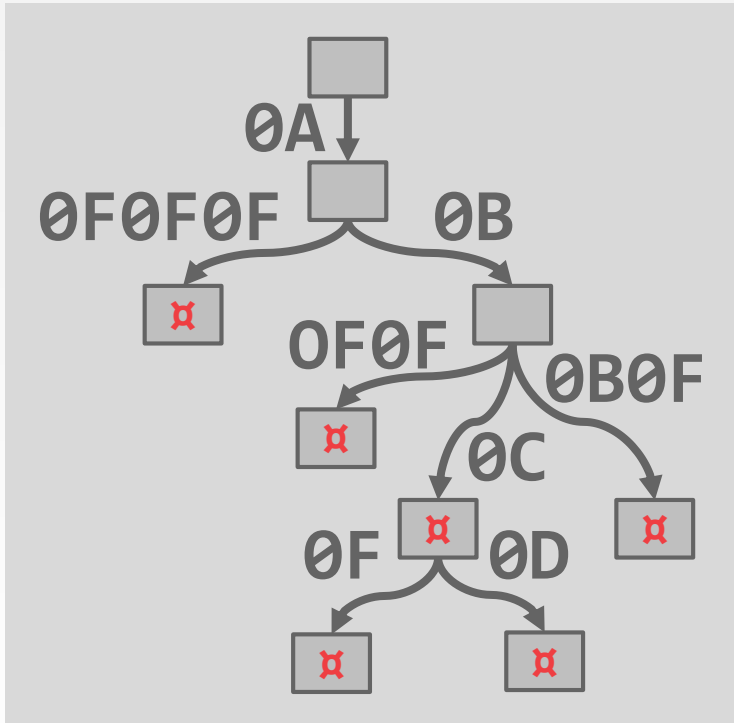Int Key: 168496141

Hex Key: 0A 0B 0C 0D

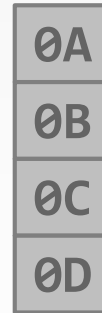| 0A |
|----|
| 0B |
| 0C |
| 0D |

*Big Endian*

| 0D |
|----|
| 0C |
| 0B |
| 0A |

*Little Endian*
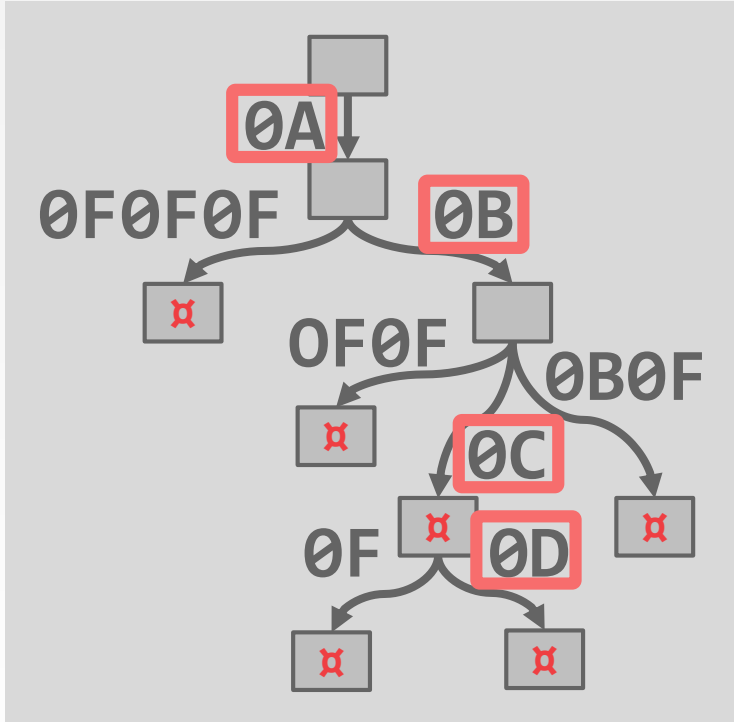
# RADIX TREE: BINARY COMPARABLE KEYS



**Big Endian**   **Little Endian**

# RADIX TREE: BINARY COMPARABLE KEYS



Int Key: 168496141

Hex Key: 0A 0B 0C 0D

| 0A |
|----|
| 0B |
| 0C |
| 0D |

*Big Endian*

| 0D |
|----|
| 0C |
| 0B |
| 0A |

*Little Endian*

CARNEGIE MELLON
**DATABASE GROUP**

# SELECTION CONDITIONS

The DBMS can use a B+Tree index if the query provides all of the attributes in a prefix of the search key.

→ Index on `<a,b,c>` matches `(a=5 AND b=3)`, but not `b=3`.

For Hash index, we must have all attributes in search key.

# B+TREE PREFIX SEARCH

*Find "XY"*

# B+TREE PREFIX SEARCH

*Find "XY"*

*Find "_Y"*

**???**

CARNEGIE MELLON
**DATABASE GROUP**

# PARTIAL INDEXES

Create an index on a subset of the entire table. This potentially reduces its size and the amount of overhead to maintain it.

```
CREATE INDEX idx_foo
          ON foo (a, b)
       WHERE c = 'WuTang'
```

CARNEGIE MELLON
DATABASE GROUP

# PARTIAL INDEXES

Create an index on a subset of the entire table. This potentially reduces its size and the amount of overhead to maintain it.

```
CREATE INDEX idx_foo
          ON foo (a, b)
       WHERE c = 'WuTang'
```

```
SELECT b FROM foo
 WHERE a = 123
   AND c = 'WuTang'
```

# COVERING INDEXES

If all of the fields needed to process the query are available in an index, then the DBMS does not need to retrieve the tuple.

```
CREATE INDEX idx_foo
          ON foo (a, b)
```

```
SELECT b FROM foo
 WHERE a = 123
```

# COVERING INDEXES

If all of the fields needed to process the query are available in an index, then the DBMS does not need to retrieve the tuple.
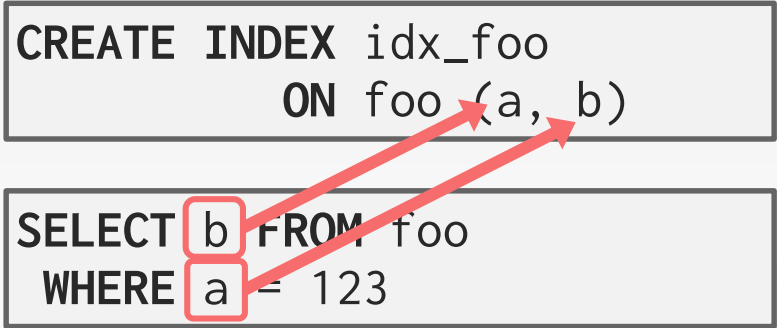
```
CREATE INDEX idx_foo
          ON foo (a, b)
```

```
SELECT b FROM foo
 WHERE a = 123
```

CARNEGIE MELLON
DATABASE GROUP

# INDEX INCLUDE COLUMNS

Embed additional columns in indexes to support index-only queries.

Not part of the search key.

```
CREATE INDEX idx_foo
         ON foo (a, b)
     INCLUDE (c)
```

```
SELECT b FROM foo
 WHERE a = 123
   AND c = 'WuTang'
```

# INDEX INCLUDE COLUMNS

Embed additional columns in indexes to support index-only queries.

Not part of the search key.

```
CREATE INDEX idx_foo
        ON foo(a, b)
    INCLUDE (c)
```
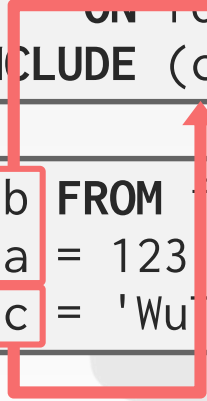
```
SELECT b FROM foo
 WHERE a = 123
   AND c = 'WuTang'
```

CARNEGIE MELLON
DATABASE GROUP

# CONCLUSION

The venerable B+Tree is always a good choice for your DBMS.

Skip Lists and Radix Trees have some interesting properties.

We will cover lock free data structures in 15-721.

# NEXT CLASS

Query Processing
→ How to use what we've talked about
so far to actually execute queries!

CARNEGIE MELLON
DATABASE GROUP