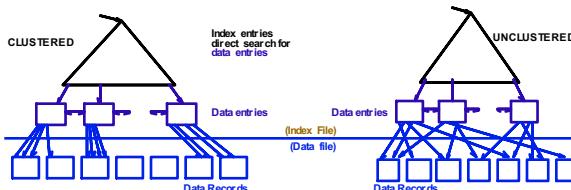


 Recall: Data Storage Considerations

- How is data stored in the index?
 - By Value:** actual data record (with key value k)
 - By Reference:** $\langle k, \text{rid of matching data record} \rangle$
 - By List of Refs.:** $\langle k, \text{list of rids of all matching data records} \rangle$
- If stored externally (alt. 2, 3)
 - CLUSTERED:** Index entries direct search for data entries
 - UNCLUSTERED:** Index entries point to Data entries (Index File) (Data file)

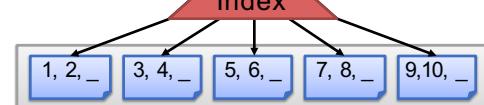


 Recall: Example from Last Lecture

Assumptions:

- Store data by reference (Alternative 2)
- Clustered tree index with 2/3 full heap file page
- Sorted** heap file
- Fan-out (F):** relatively large.
- Assume static index

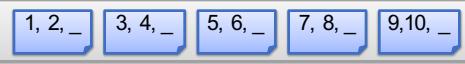
How do we build this?



 Simple Idea?

Input Heap File 

- Sort heap file & leave some space

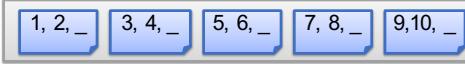


- Pages in logical order (seq. disk access)
- Do we need next pointers?
 - Pages are in logical order (might still add them for simplicity)

 Simple Idea?

Input Heap File 

- Sort heap file & leave some space

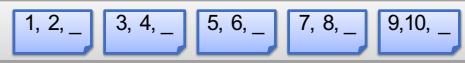


- Cost?** (Assume $> \sqrt{\# \text{Pages}}$ memory)
 - How many passes?
 - Two passes

 Simple Idea?

Input Heap File 

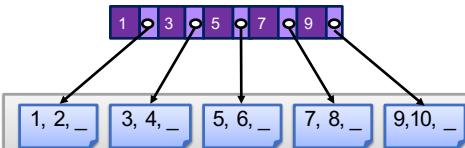
- Sort heap file & leave some space



- Step 2: Build the index...**
 - Why not just binary search heap file?
 - Load record data just to read key
 - Small fan-out \rightarrow deep tree \rightarrow more IOs

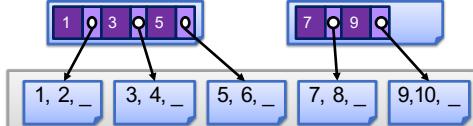
 Build a high fan-out search tree

- Start simple: Sorted (key, page id) file
 - No record data
 - Binary search key file
 - Forgot:** Need to break across pages!



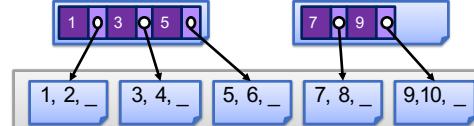
 Build a high fan-out search tree

- Start simple: Sorted (key, page id) file
 - No record data
 - Binary search key file
 - Forgot:** Need to break across pages!



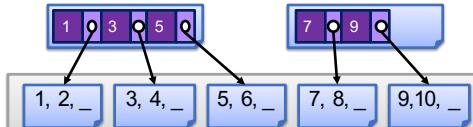
 Build a high fan-out search tree

- Start simple: Sorted (key, page id) file
 - No record data
 - Binary search key file
 - Complexity?**



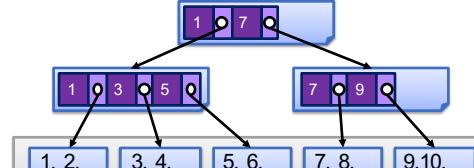
 Build a high fan-out search tree

- Start simple: Sorted (key, page id) file
 - No record data
 - Binary search key file
 - Complexity:** Still binary search of #Pages / (#Pair<key,ptr> that fit on page)



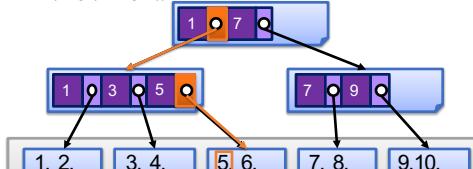
 Build a high fan-out search tree

- Recursively “index” key file
- Key Invariant: \leftarrow Pun intended
 - Node $[..., (K_L, P_L), (K_R, P_R), ...]$ \rightarrow All tuples in range $K_L \leq K < K_R$ are in tree P_L



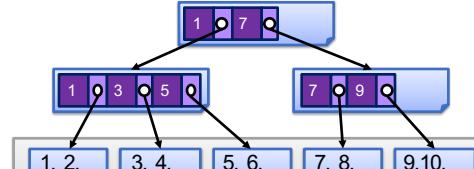
 Search a high fan-out search tree

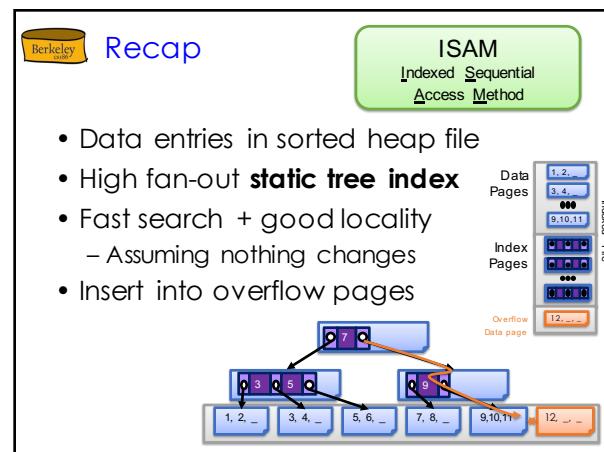
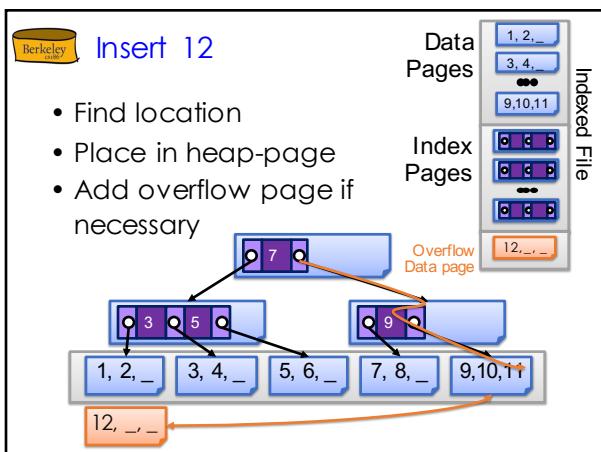
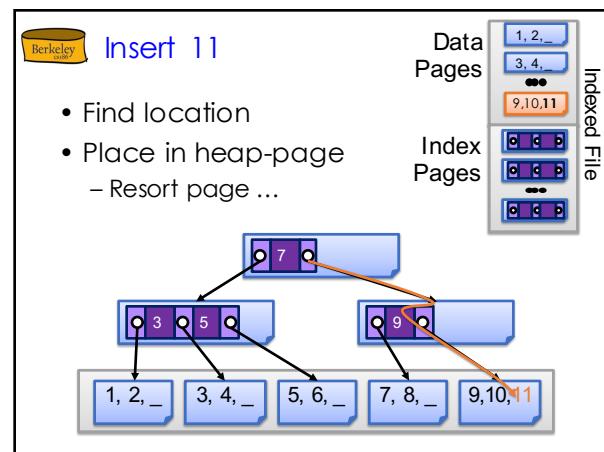
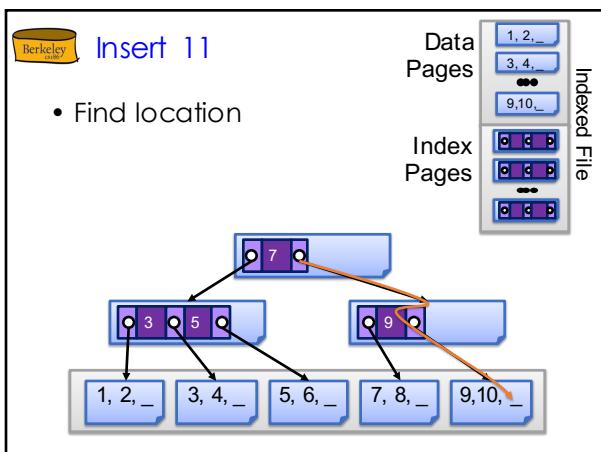
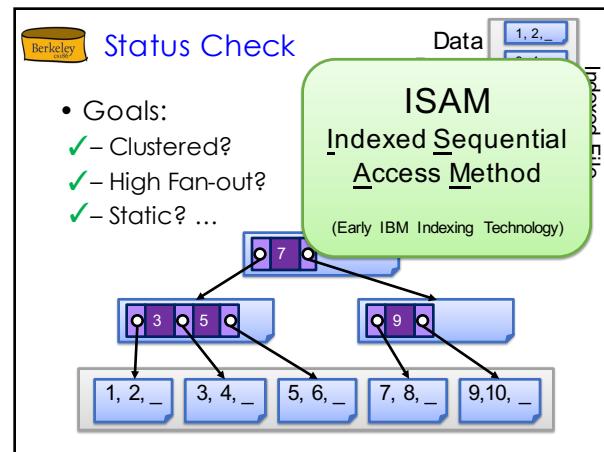
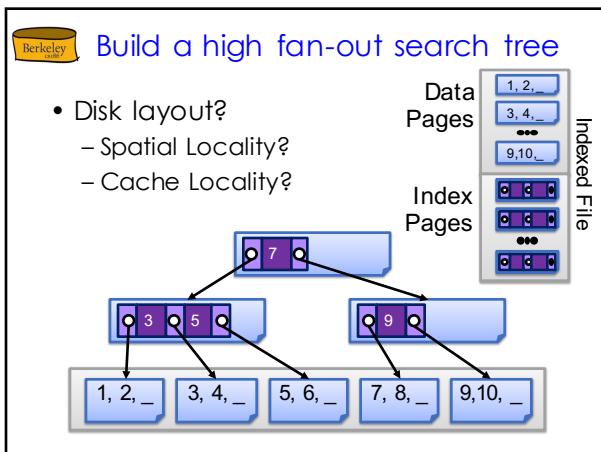
- Searching for 5?
 - Binary Search each node (page) starting at root
 - Follow pointers to next level of search tree
- Complexity?**
 - $O(\log_2(\#Pages))$



 Left Key Optimization

- Optimization:
 - Do we need the left most key?

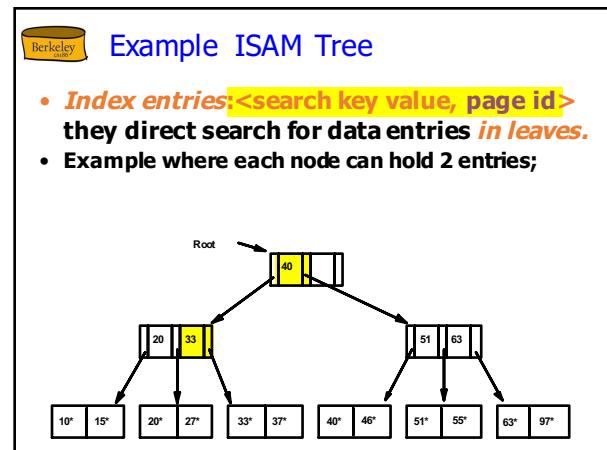
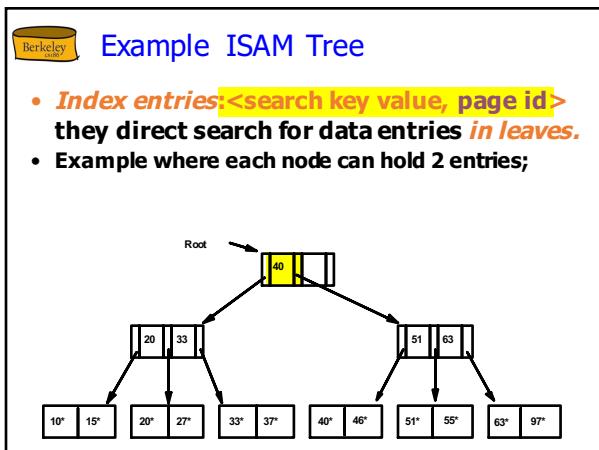
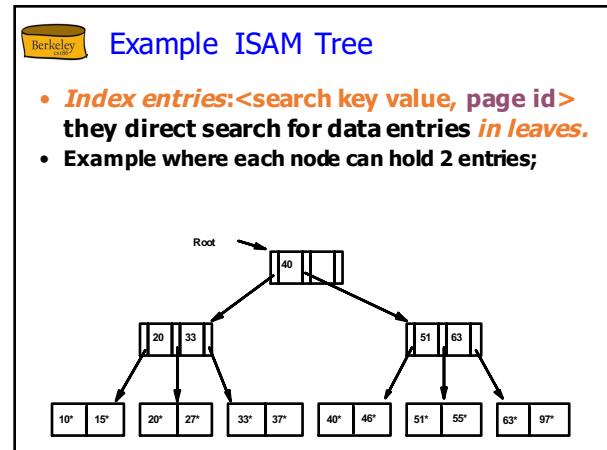
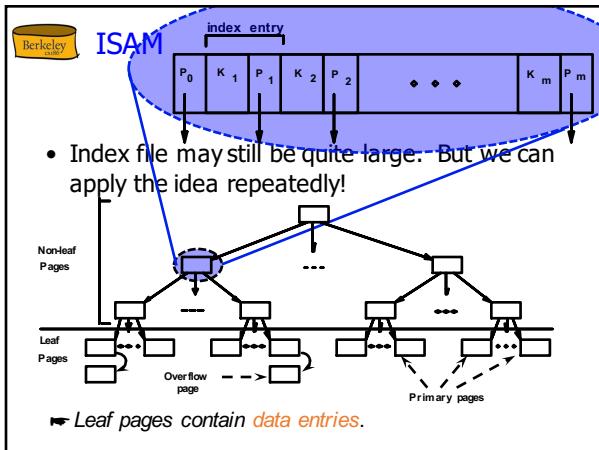
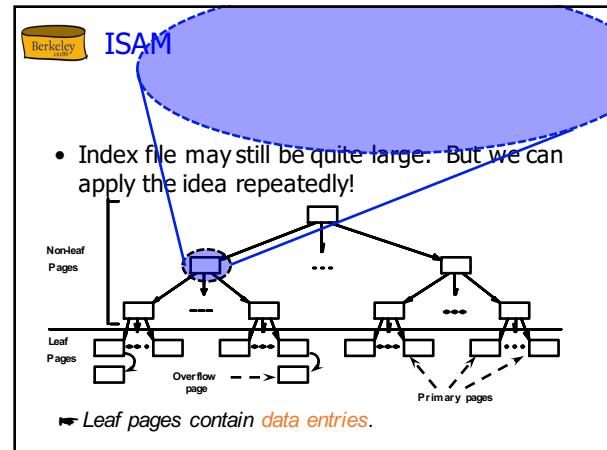


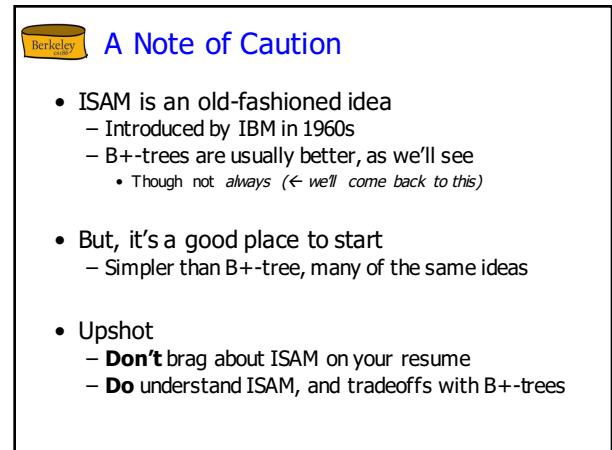
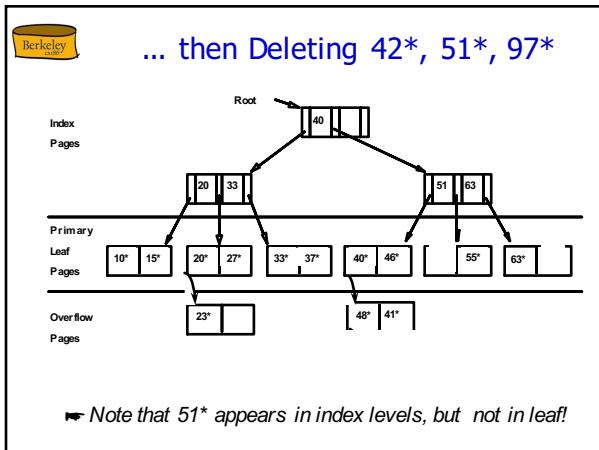
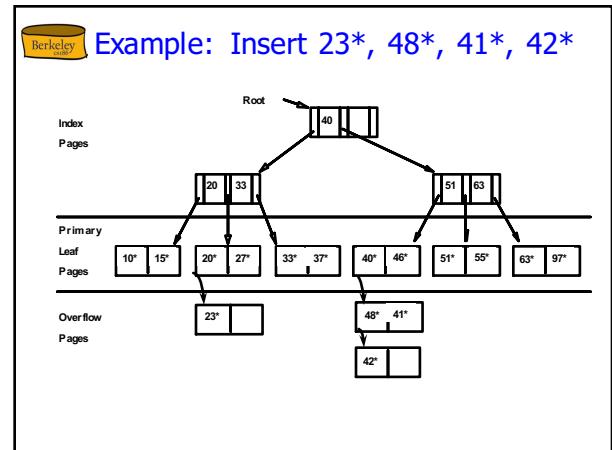
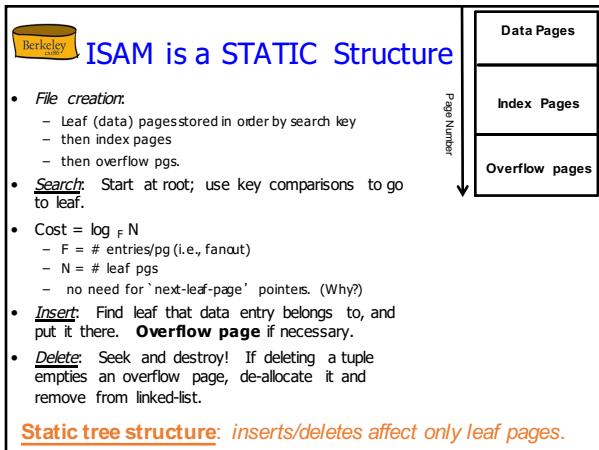
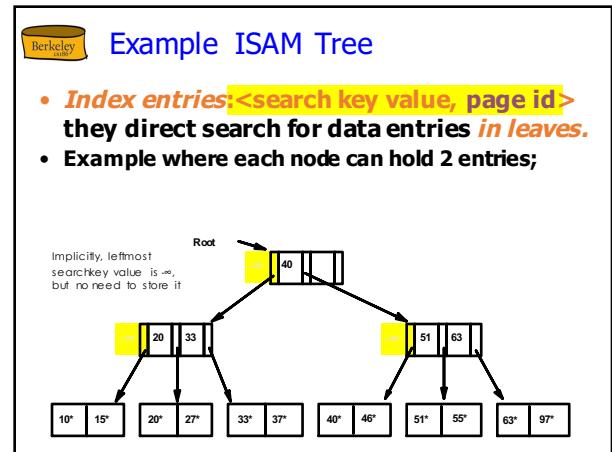
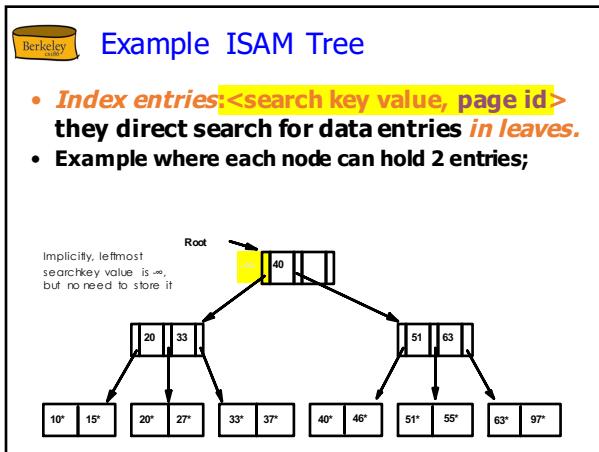


Range Searches

- “Find all students with gpa > 3.0”
 - Sorted file? Binary search to find first, scan to find others.
 - Cost of binary search in a database can be quite high. Q: Why?
- Simple idea: Create an ‘index’ file.

Can do binary search on (smaller) index file!

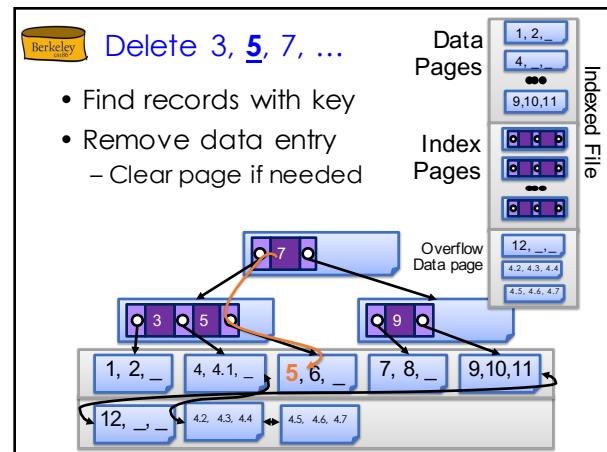
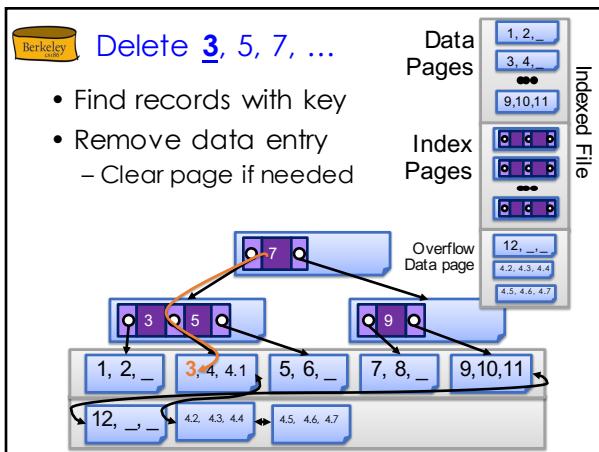
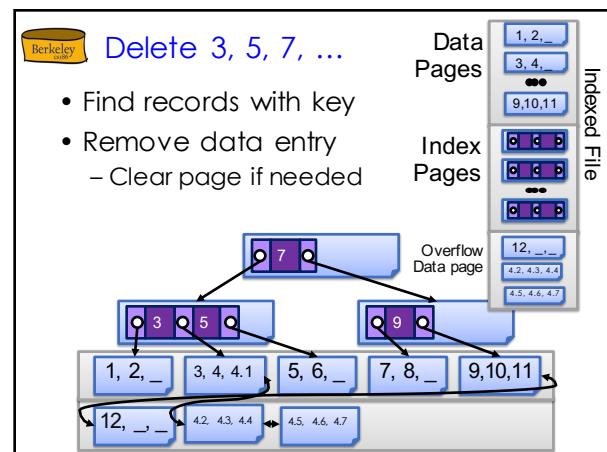
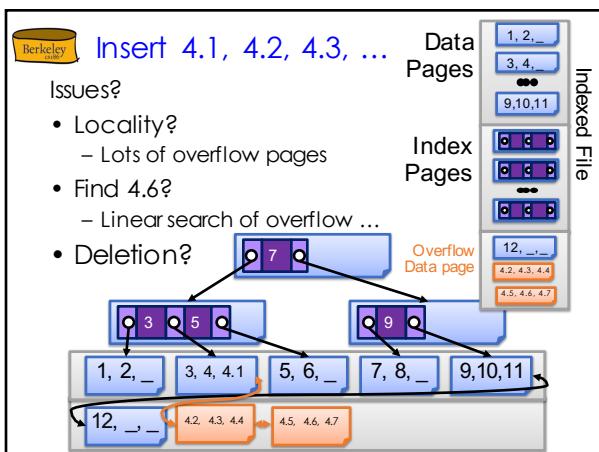
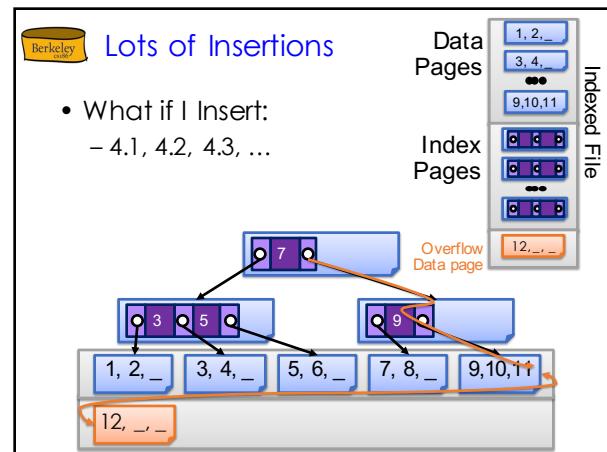


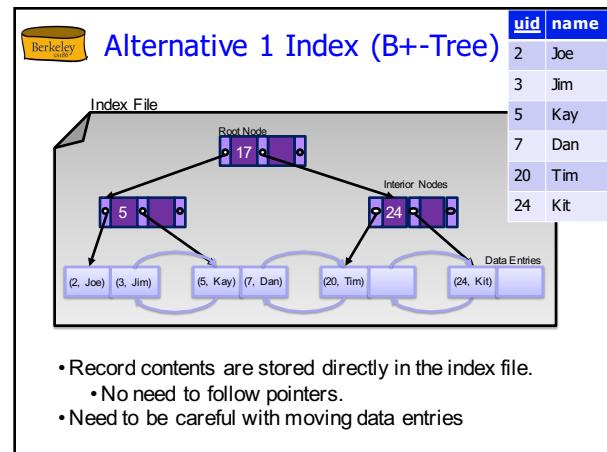
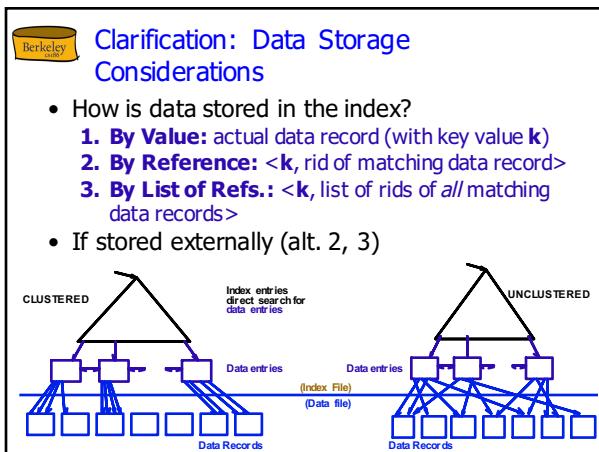
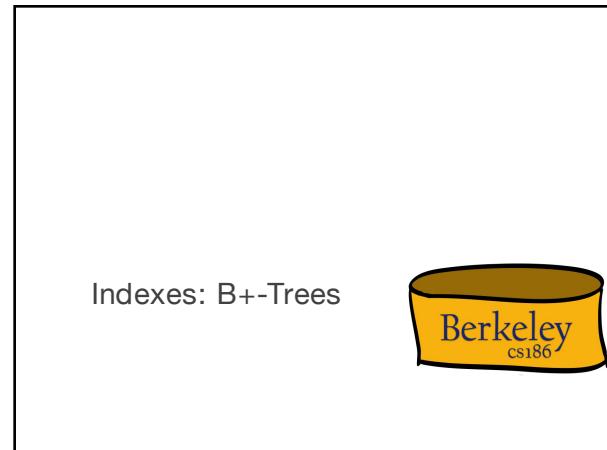
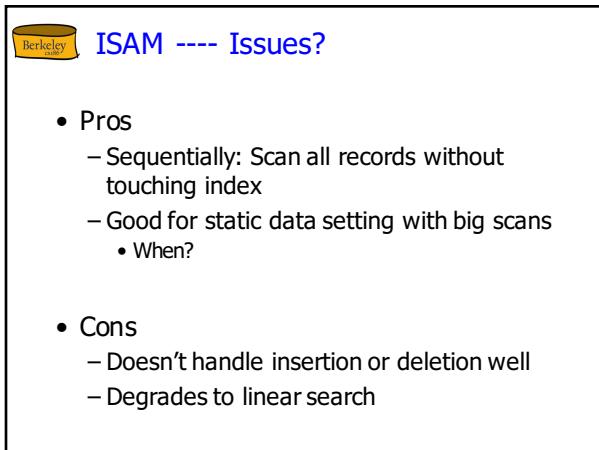
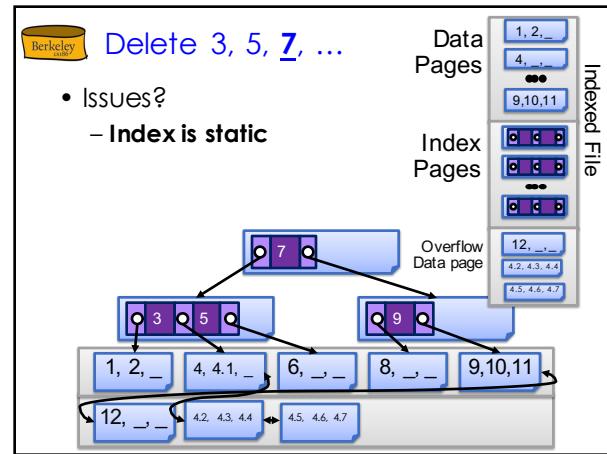
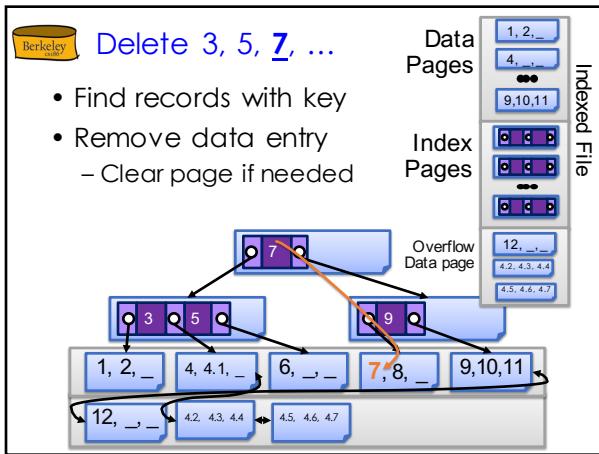


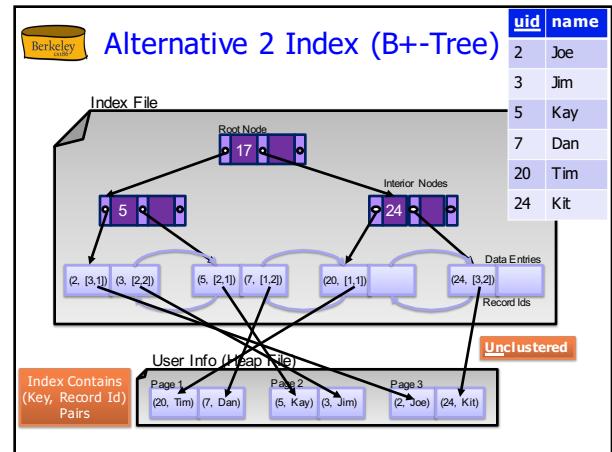
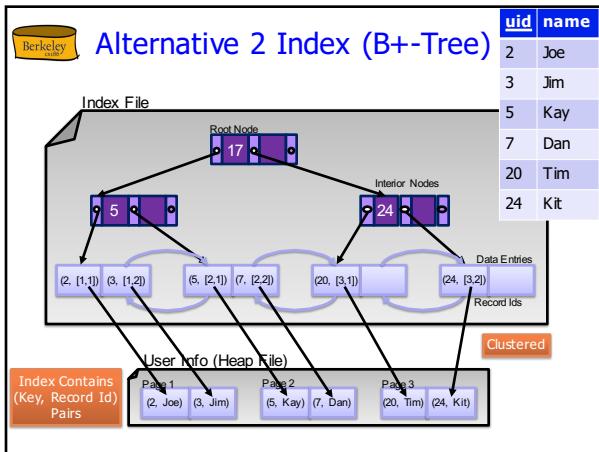
A Note of Caution

- ISAM is an old-fashioned idea
 - Introduced by IBM in 1960s
 - B+-trees are usually better, as we'll see
 - Though ...
- But, it's a ...
 - Simpler
- Upshot
 - **Don't** brag about ISAM on your resume
 - **Do** understand ISAM, and tradeoffs with B+-trees

What's wrong with ISAM?







Berkeley **Enter the B+-Tree**

- Similar to ISAM
 - Same interior node structure:
 - <Key, Page Ptr> Pairs with same guarantees
 - Same search routine as before
- Dynamic Tree Index**
 - Always Balanced
 - Support efficient insertion deletion
 - Grows at root not leaves!
- “+”? B-tree that stores data in leaves only
 - Recall: Data = Records, Ref, or List of Ref

Berkeley **What does B stand for?**

- Introduced in 1970 by Rudolf Bayer and Edward M. McCreight
- <https://vimeo.com/73481096>
 - Scrub to 16 minutes

Berkeley **Example of a B+-Tree**

- Structure is similar to ISAM
- Occupancy Invariant:**
Each interior node is at least partly full:
 $d \leq \# \text{ entries} \leq 2d$
 - d : **order of the tree** (max fan-out = $2d + 1$)
- Data pages at bottom may not be in logical
 - Next and last pointers required

Root Node

Data Pages

Height 1: 5 x 4 = 20 Records

Berkeley **B+ Trees and Scale**

- How big is a height 2 B+-Tree
 - $d = 2 \rightarrow$ Fan-out?
 - Fan-out = $2d + 1 = 5$

Root Node

Height 1: 5 x 4 = 20 Records

B+ Trees and Scale

- How big is a height 2 B+-Tree
 - $d = 2 \rightarrow$ Fan-out?
 - Fan-out = $2d + 1 = 5$

B+ Trees in Practice

- Typical order: 1600. Typical fill-factor: 67%.
 - average fan-out = 2144
 - (assuming 128 Kbytes pages at 40Byte per record)
- At typical capacities:
 - Height 1: $2144^2 = 4,596,736$ records
 - Height 2: $2144^3 = 9,855,401,984$ records
- Can often hold top levels in buffer pool:
 - Level 1 = 1 page = 128 Kbytes
 - Level 2 = 2144 pages = 274.4 Mbyte
 - Addresses 588.38 Gigabytes of leaf (data entry) pages.

Searching the B+-Tree

- Same as ISAM
- Find key = 27
 - Find split on each node
 - Follow pointer to next node

Use binary search on each page

Am I done?
(Alternative 1 vs 2,3)

Searching the B+-Tree

Searching the B+-Tree

Table (Heap File)

Page 1 (20, Tim)	Page 2 (7, Dan)	Page 3 (5, Kay)	Page 4 (3, Jim)
(27, Joe)	(34, Kit)	(1, Kim)	(42, Hal)

Alternative 2
Recorded

Searching the B+-Tree

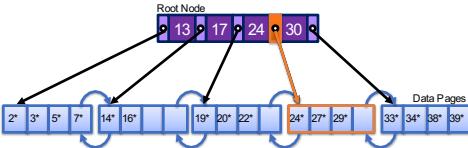
- Same as ISAM
- Find key = 27
 - Find split on each node
 - Follow pointer to next node

Use binary search on each page

How about insert?
(are you ready?)

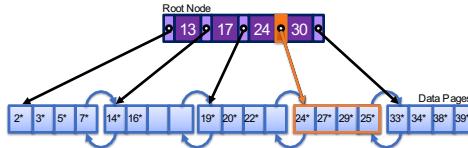
 **Inserting 25* into a B+-Tree**

- Find the correct leaf



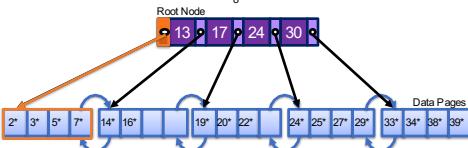
 **Inserting 25* into a B+-Tree**

- Find the correct leaf
- If there is room in the leaf just add the entry
 - Sort the page leaf page by key



 **Inserting 8* into a B+-Tree**

- Find the correct leaf

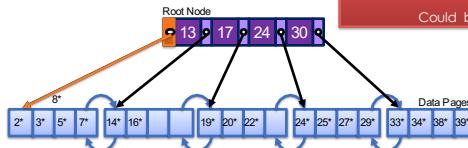


 **Inserting 8* into a B+-Tree**

- Find the correct leaf
 - Split leaf if there is not enough room
 - Redistribute entries evenly

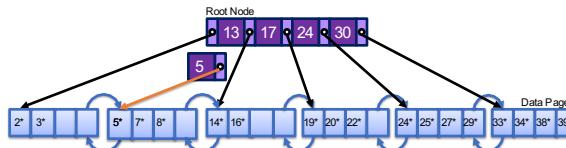
What happens to record ids if this is an Alternative 1 index (store data by value)?

Careful:
Moving records to new pages changes records id
Could be pricey!!!



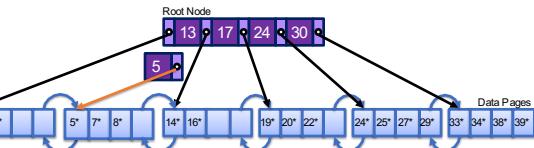
 **Inserting 8* into a B+-Tree**

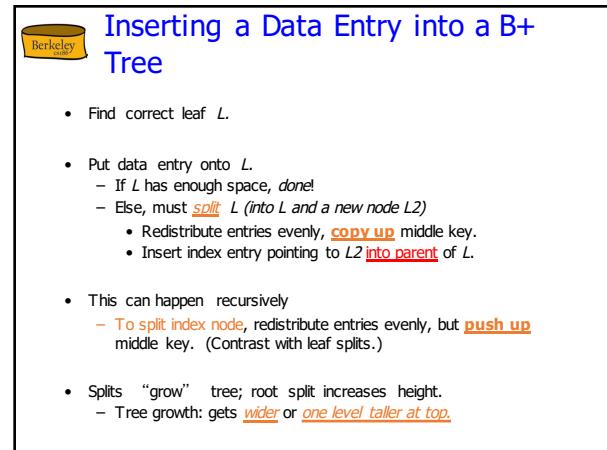
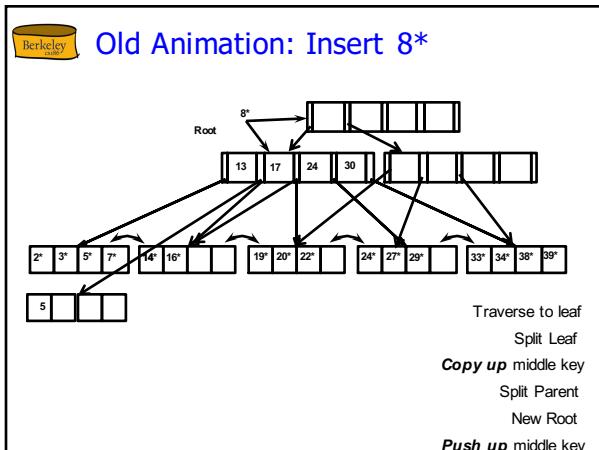
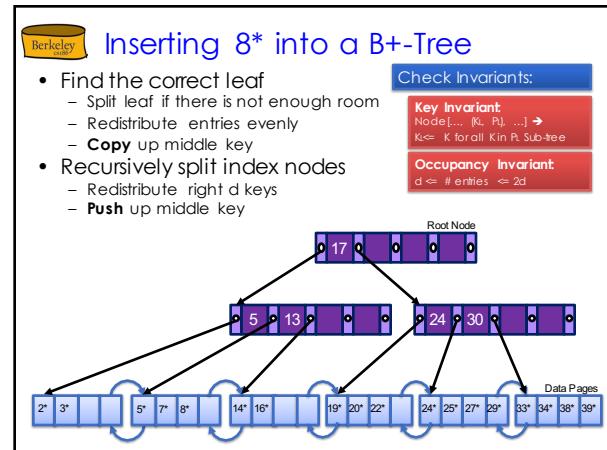
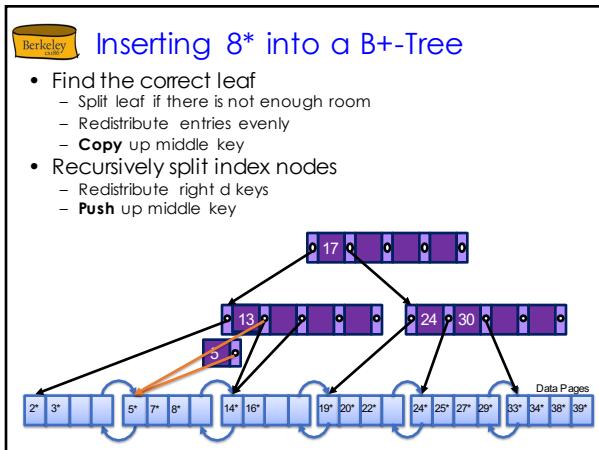
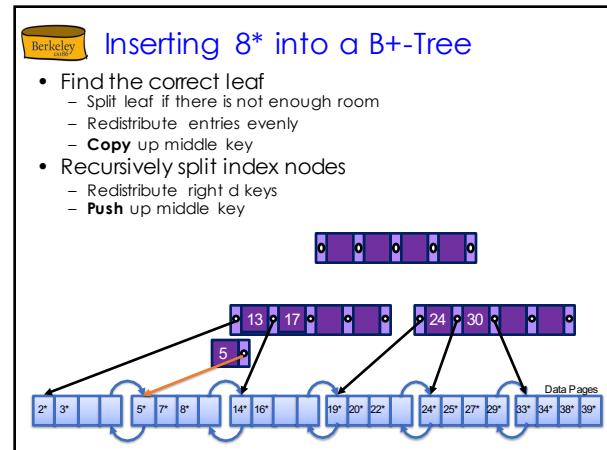
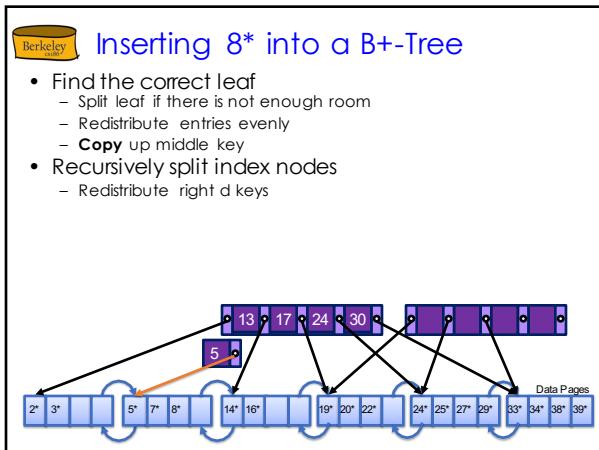
- Find the correct leaf
 - Split leaf if there is not enough room
 - Redistribute entries evenly
 - Copy** up middle key

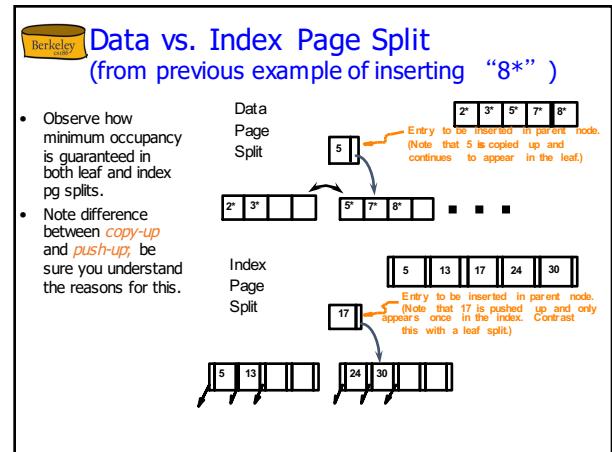
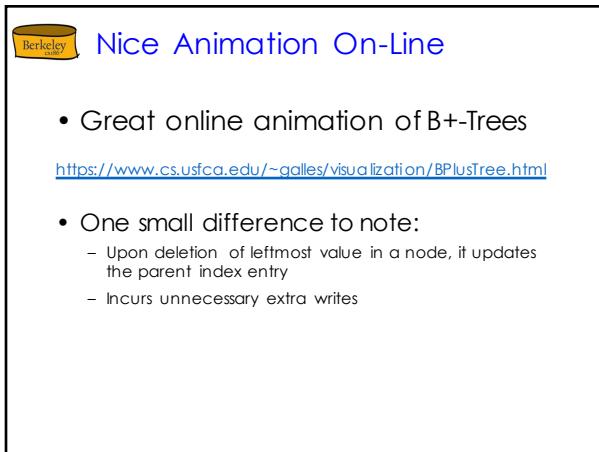
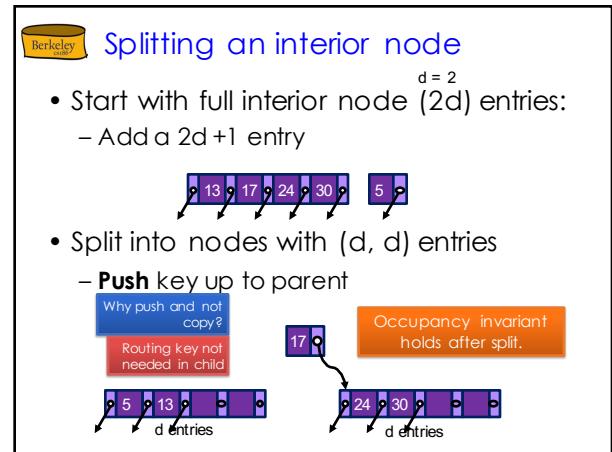
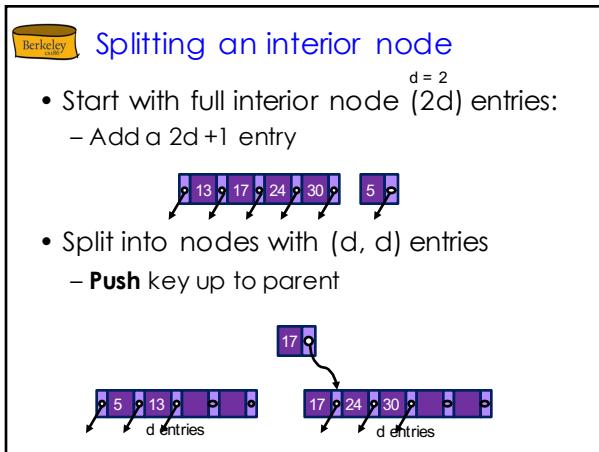
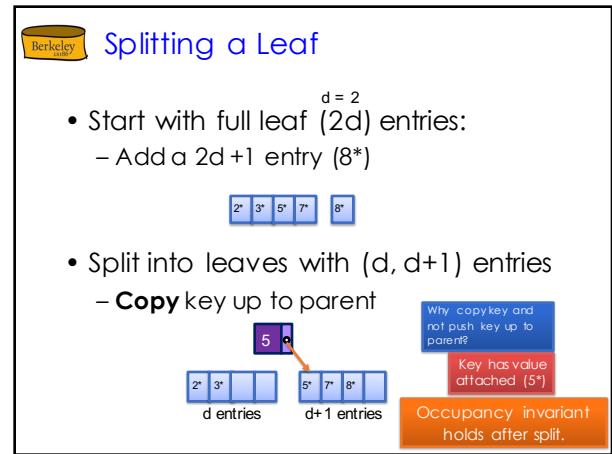
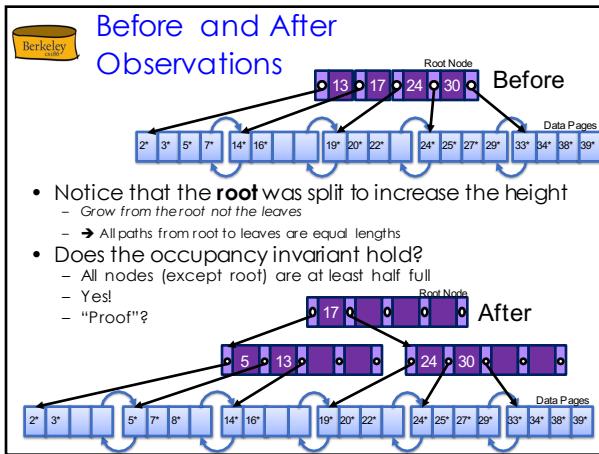


 **Inserting 8* into a B+-Tree**

- Find the correct leaf
 - Split leaf if there is not enough room
 - Redistribute entries evenly
 - Copy** up middle key

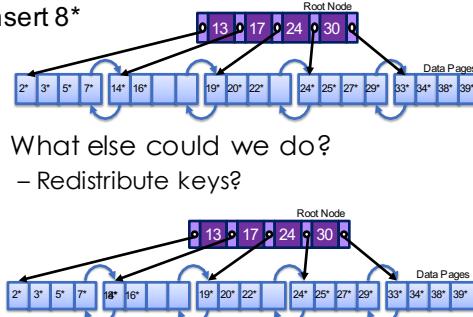






 Did we have to split?

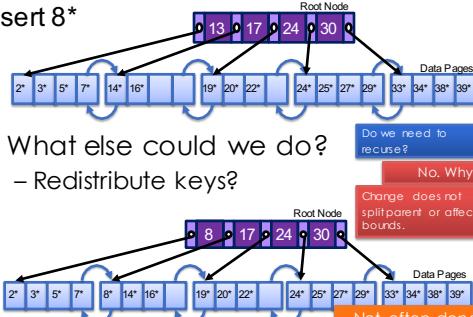
Insert 8*



- What else could we do?
 - Redistribute keys?

 Did we have to split?

Insert 8*



- What else could we do?
 - Redistribute keys?

Do we need to recurse?
No. Why?
Change does not split parent or affect bounds.

Not often done in practice

 Bulk Loading of B+ Tree

Suppose we want to build an index on a large table

- Would it be efficient to just call insert repeatedly?
 - No ... Why not?
 - Random Order: CLZARNDXEKFWIUB

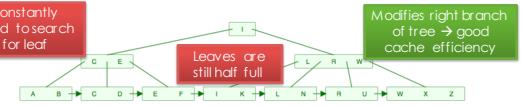


Work through this animation <https://www.cs.usc.edu/~gales/visualization/BPlusTree.htm>

 Bulk Loading of B+ Tree

Suppose we want to build an index on a large table

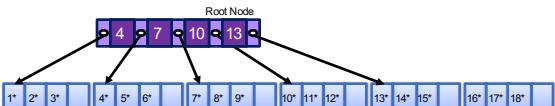
- Would it be efficient to just call insert repeatedly?
 - No ... Why not?
 - **Sorted** Order? ABCDEFKLNRUWXZ



Work through this animation <https://www.cs.usc.edu/~gales/visualization/BPlusTree.htm>

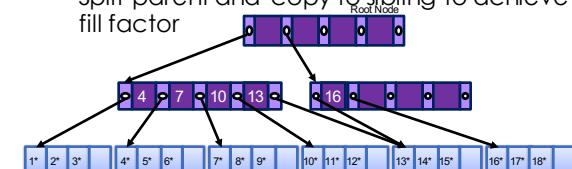
 Smarter Bulk Loading a B+-Tree

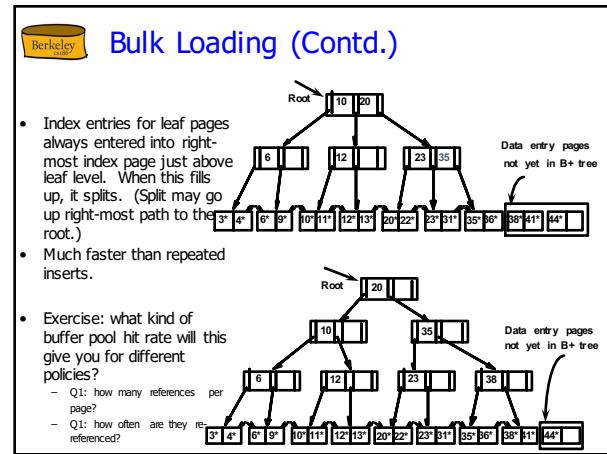
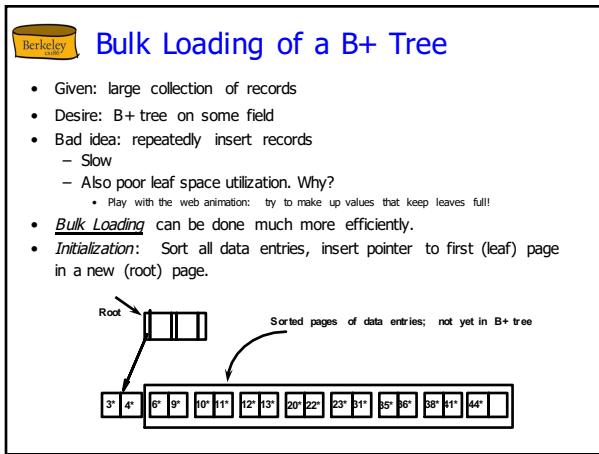
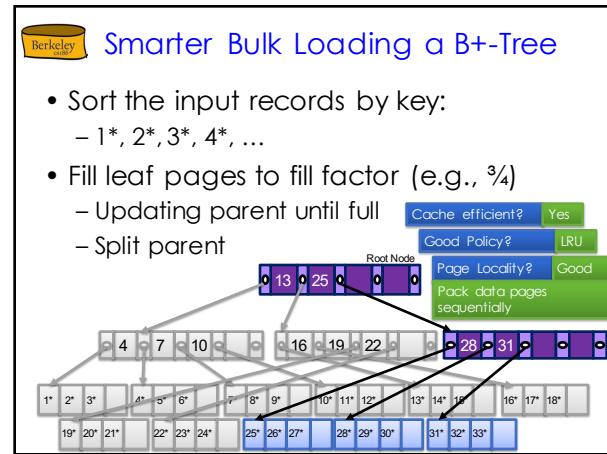
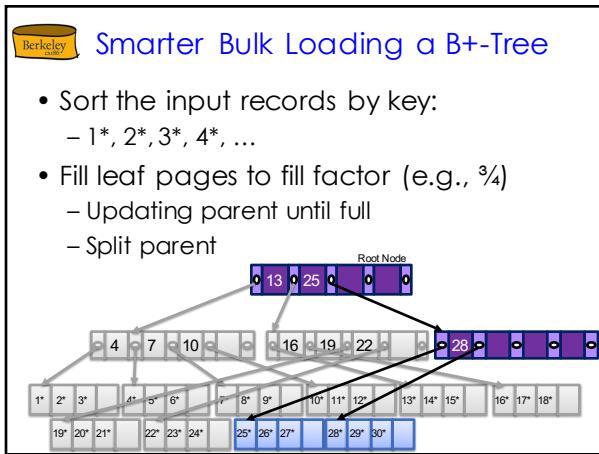
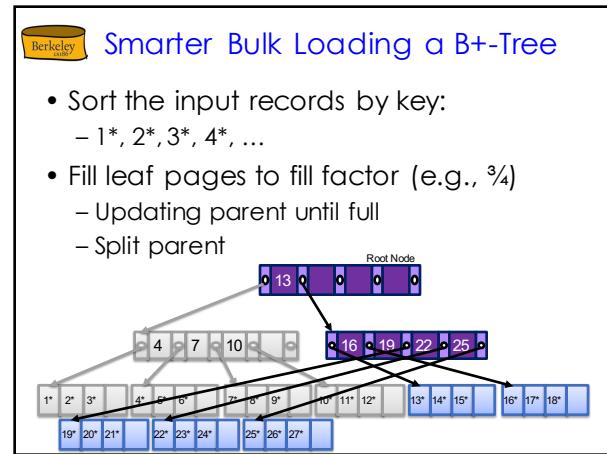
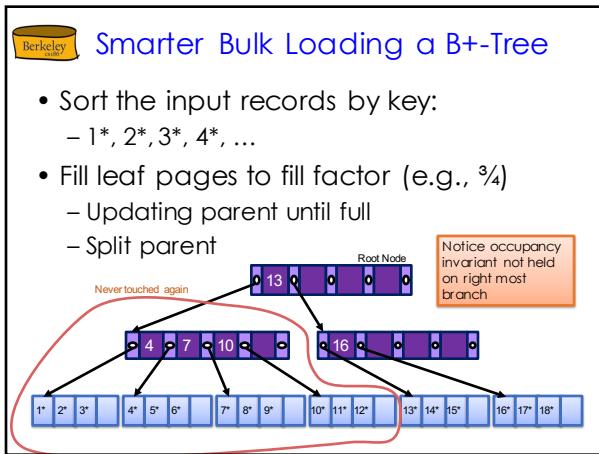
- Sort the input records by key:
 - 1*, 2*, 3*, 4*, ...
- Fill leaf pages to fill factor (e.g., $\frac{3}{4}$)
 - Updating parent until full



 Smarter Bulk Loading a B+-Tree

- Sort the input records by key:
 - 1*, 2*, 3*, 4*, ...
- Fill leaf pages to fill factor (e.g., $\frac{3}{4}$)
 - Updating parent until full
 - Split parent and copy to sibling to achieve fill factor





 **Summary of Bulk Loading**

- Option 1: multiple inserts.
 - Slow.
 - Does not give sequential storage of leaves.
- Option 2: Bulk Loading
 - Fewer IOs during build. (Why?)
 - Leaves will be stored sequentially (and linked, of course).
 - Can control “fill factor” on pages.

 **Deleting a Data Entry from a B+ Tree**

- Start at root, find leaf L where entry belongs.
- Remove the entry.
 - If L is at least half-full, done!
 - If L has only $d-1$ 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 .
- Merge could propagate to root, decreasing height.

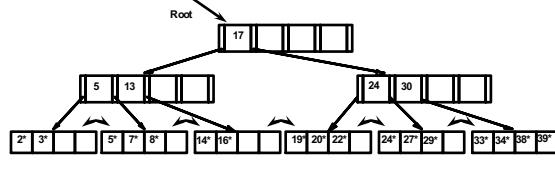
 **Deleting a Data Entry from a B+ Tree**

- Start at root and leaf L where entry belongs.
- Remove the entry.
 - If L is at least half-full, done!
 - If L has only $d-1$ 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 .
- Merge could propagate to root, decreasing height.

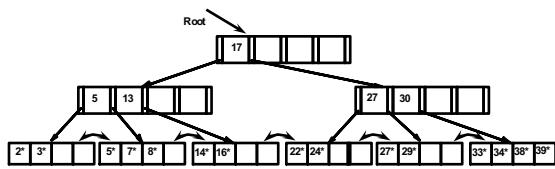
In practice, many systems do not worry about ensuring half-full pages. Just let page slowly go empty; if it's truly empty, just delete from tree and leave unbalanced.

You won't be tested on delete.

 **Example Tree (including 8*) Delete 19* and 20* ...**



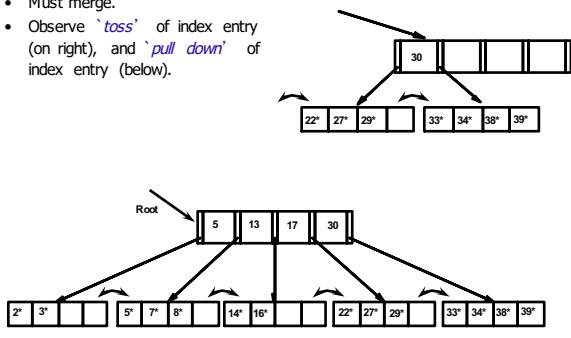
 **Example Tree (including 8*) Delete 19* and 20* ...**



- Deleting 19* is easy.
- Deleting 20* is done with re-distribution. Notice how middle key is copied up.

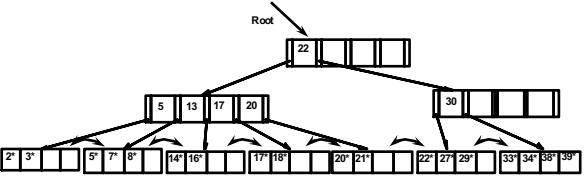
 **... And Then Deleting 24***

- Must merge.
- Observe ‘toss’ of index entry (on right), and ‘pull down’ of index entry (below).



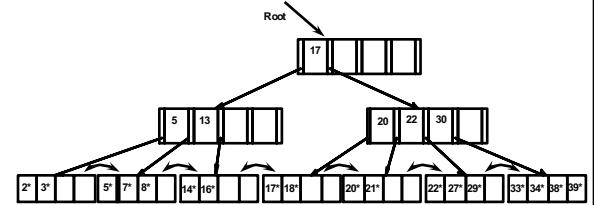
 Example of Non-leaf Re-distribution

- Tree is shown below *during deletion* of 24*. (What could be a possible initial tree?)
- In contrast to previous example, can re-distribute entry from left child of root to right child.



 After Re-distribution

- Intuitively, entries are **re-distributed by ‘pushing through’** the splitting entry in the parent node.
- It suffices to re-distribute index entry with key 20; we’ve re-distributed 17 as well for illustration.



 Variable Length Keys & Records

- So far we have been using integer keys



- What would happen to our **occupancy** invariant with variable length keys?

Dan Ha | Dannon Yogurt | Davey Jones | David Yu | Devarakonda Murthy

- What about data in leaf pages:

Dan Ha : {3, 14, 30, 50, 75, 90}	Dan He: {12}	Dan Ham: {1}	Dannon Smith: {}
----------------------------------	--------------	--------------	------------------

 Redefine Occupancy Invariant

- Order (d)** makes little sense with variable-length entries
 - Variable sized** records and search keys:
 - different nodes have different numbers of entries.
 - Example:** single entry might fill entire page
 - Index pages** often hold many **more entries** than leaf pages.
 - Even with fixed length fields, Alternative (3) gives variable length data entries
- Use a physical criterion in practice: **at least half-full**
 - Measured in **bytes**
- Many real systems are even sloppier than this --- only reclaim space when a page is **completely** empty.
 - Why?
 - Reduce unnecessary IOs

 Prefix Compress Keys?

- How can we get more keys on a page?

Dan Ha | Dannon Yogurt | Davey Jones | David Smith | Devarakonda Murthy

- What if we compress the keys:



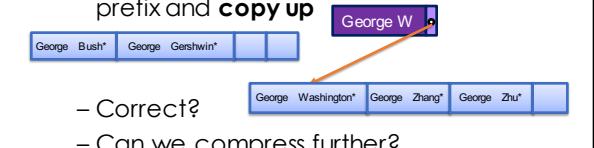
- Are these the same?
 - David Jones?
- Not the same! Can't just reduce to prefix.
- But maybe we can fix it?

 Prefix Key Compression

- What if we compress starting at leaf:

George Bush*	George Gershwin*	George Washington*	George Zhang*	George Zhu*
--------------	------------------	--------------------	---------------	-------------

- on split, determine minimum splitting prefix and **copy up**



- Correct?
- Can we compress further?

George A	George B	George S	George Sm	George W
----------	----------	----------	-----------	----------

 **Prefix Key Compression**

- Important to increase fan-out. (Why?)
- Key values in index entries just 'direct traffic'; can often compress them.

Dannon Yogurt	Davey Jones	David Smith	Devarakonda Murthy
---------------	-------------	-------------	--------------------

 **Prefix Key Compression**

- Important to increase fan-out. (Why?)
- Key values in index entries just 'direct traffic'; can often compress them.

Dan	Dave	Davi	De	
-----	------	------	----	--

- Is this correct?
- Ensure that each index entry is greater than every key value (in any descendant leaf) to its left.
- In practice we compress upon "copy up" from leaf.

 **Suffix Key Compression**

- All keys have large common prefix

George A	George B	George S	George Sm	George W
----------	----------	----------	-----------	----------

- Move common prefix to header

"George "	A	B	S	Sm	W				
-----------	---	---	---	----	---	--	--	--	--

- Still use full prefix in comparisons
- When might this be especially useful?
 - Composite Keys. Example?
 - <Zip code, Last Name, First Name>

 **Suffix Key Compression**

- If many index entries share a common prefix

MacDonald	MacDougal	MacFeeley	MacLaren	
-----------	-----------	-----------	----------	--

- Particularly useful for composite keys
 - Why?

 **Suffix Key Compression**

- If many index entries share a common prefix

Mac	Donald	Dougal	Feeley	Laren	
-----	--------	--------	--------	-------	--

- Particularly useful for composite keys
 - Why?

 **Summary**

- ISAM is a static structure.
 - Only leaf pages modified**; overflow pages needed.
 - Overflow **chains can degrade performance** unless size of data set and data distribution stay constant.
- B+ tree is a dynamic structure.
 - Inserts/deletes leave tree height-balanced; $\log_F N$ cost.
 - High fanout (F) means depth rarely more than 3 or 4.
 - Almost always better than maintaining a sorted file.
 - Typically, 67% occupancy on average.
 - Usually preferable to ISAM; adjusts to growth gracefully.
 - Caution: if data entries are data records, splits can change rids!



Summary (Contd.)

- Bulk loading can be much faster than repeated inserts for creating a B+ tree on a large data set.
- Key compression increases fanout, reduces height.
- B+ tree widely used because of its versatility.
 - One of the most optimized components of a DBMS.



Architecture of a DBMS

This lecture:

- Finished **tree indexes**: ISAM and B+-Trees

Next Lecture:

- Relational Operators

