

= B-Tree = Basics, Insert

Think of "B" as "branching", not binary (unofficial)

Each node has k children ("k-ary tree") / page

Goal: Each node completely occupies one sector/unit on a hard drive: MINIMIZE # OF DISK ACESSES

↳ eval running time based on 2 things: compute time
disk accesses

Def: B-Tree (see picture on page #4)

∃ a root node - all nodes are attached (directly or indirectly) to it

Each node has: (#) \rightarrow AND data-

◦ int cKeys - count of keys in this node

◦ int [] rgKeys - array of keys, sorted into (ascending) order \leftarrow length = $2t+1$

◦ boolean fLeaf - true if this is a leaf node, false otherwise

Nonleaf nodes have ($c_{\text{Keys}} + 1$) references to child nodes - BTNode[] $\xrightarrow{\text{rgChildren}}$

Leaf nodes have no children - rgChildren isn't defined

• EVERY LEAF HAS THE SAME DEPTH (=height of tree, " h ")

Each node has a minimum degree t (\forall nodes have same degree)

◦ Every node (EXCEPT the root) must have at least $t-1$ keys

◦ " " can have NO MORE than $2t-1$ keys

(node w/ $2t-1$ keys is FULL \Leftrightarrow has $2t$ children)

Note: height of B-Tree is $\log_t(N)$

DiskPages read:

$O(h)$ =

$O(\log_t N)$

ComputeTime = $O(t)$

$O(t \cdot \log_t N)$

=BTTreeSearch (BTreeNode x, int k) $\xrightarrow{\text{target}}$

{ for(int i = 0; i < x.cKeys && k > rgKeys[i]; i++)

\uparrow // empty on purpose - all work done in ↑

 if(i < x.cKeys && k == rgKeys[i])

 return x.rgKeys[i].data; //awk. syntax :)

 if(x.fLeaf) return null;

 else return BTTreeSearch(DISK-READ(x.rgChildren[i]), k)

\downarrow

2/4

==== BTreeCreate () =====

```
{ X = new BTreeNode();  
  // X.fLeaf = true, in constructor  
  // X.cKeys = 0; // also "alloc" arrays  
  DISK-CREATE(X);  
  root = X;
```

// int t; global variable + contains degree

===== BTreeSplit Child (BTreeNode parent, int iChild, BTreeNode child) =====

// PRE: parent is a NON FULL(internal) node

// child is a FULL child of parent (may or may not be a leaf)

// iChild index (in parent) of child

// POST: child has been split into 2 nodes, median value stored in parent.

BTreeNode next = new BTreeNode()

next.fLeaf = child.fLeaf; next.cKeys = t - 1;

```
[ for(int i=0; i<(t-1); i++) next.rgKeys[i] = child.rgKeys[i+(t-1)];  
  if(!child.fLeaf) // internal, non leaf node - copy child refs.  
    for(int i=0; i<(t-1); i++) next.rgChildren[i] = child.rgChildren[i+(t-1)];
```

child.cKeys = t - 1;

for(int i=parent.cKeys; i>iChild; i--) // move everything up one...

```
{ parent.rgKeys[i] = parent.rgKeys[i-1];
```

parent.rgChildren[i+1] = parent.rgChildren[i-1];

}

parent.rgChildren[iChild+1] = next; // add next to tree

parent.rgKeys[iChild] = child.rgKeys[t - 1];

parent.cKeys++;

DISK-CWRITE(parent); DISK-CWRITE(child); DISK-WRITE(next);

Copy top 1/2
into new node

(internal
addition)

3 / 4

BTree Insert (BTree tree, int key)

{ // PRE: + is a valid B-Tree ; key is the value to insert
(may be an object/char.)

// POST: + is still a valid BTee; key has been inserted

// This is a 'one pass' implementation - if we notice a child is
// full, we will preemptive split it, BEFORE traversing it

```
BTreeNode *t = tree.root;
if( root.cKeys == 2*t-1) // need to split root
{
    BTreeNode newRoot = new BTreeNode();
    newRoot.fLeaf = false;
    newRoot.cKeys = 0;
    newRoot.Children[0] = root;
    BTreeSplitChild( newRoot, 0, root );
    tree.root = newRoot;
    root = newRoot;
}
```

3 ← BTreeInsertNonfull(root, 1);

BTree Insert Nonfull (BTreeNode node, int key)

```
{ if( node.fLeaf == true )
    {
        for( int i = node.cKeys; i >= 1 && key < node.rgKeys[i-1]; i-- )
            node.rgKeys[i] = node.rgKeys[i-1];
        node.rgKeys[i] = key;
        node.cKeys++;
        DISK-WRITE( node );
    }
}
```

3 else // internal node - figure out which child to add the key to

```
for ( int i = node.cKeys-1; i >= 0 && key < node.rgKeys[i]; i-- )
```

^{/* for loop empty */} i++ ^{// move i up 1 after loop}

```
BTreeNode child = DISK-READ( node.rgChildren[i] );
```

if (child.cKeys == 2*t-1)

```
BTreeSplitChild( node, i, child )
```

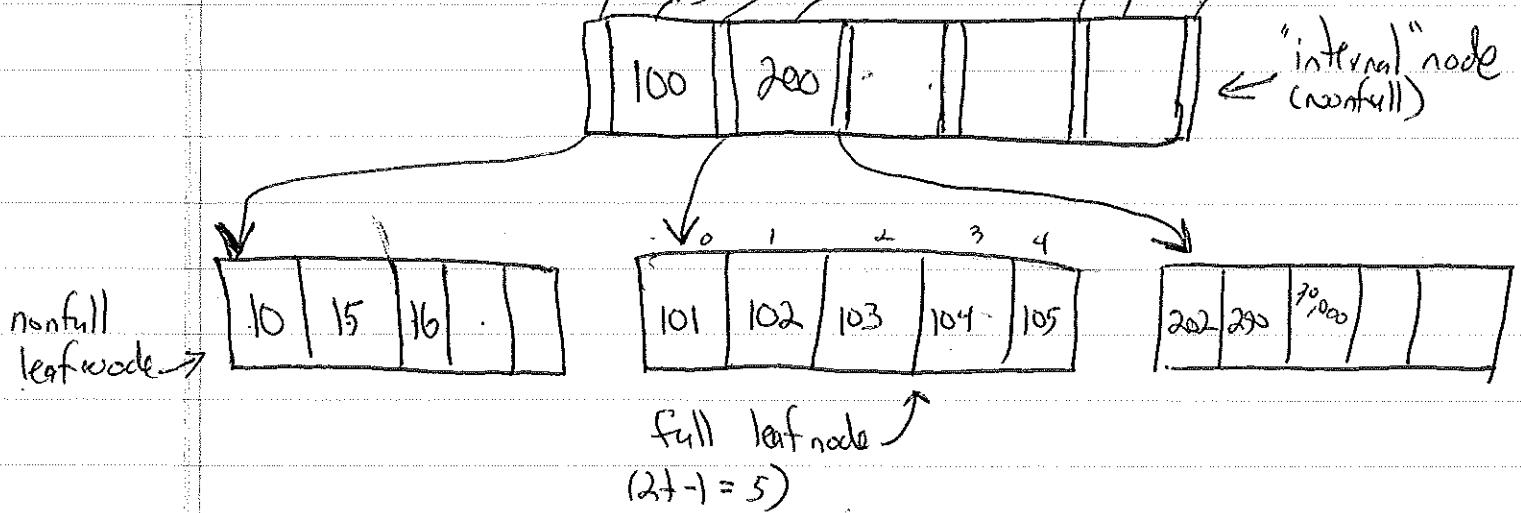
if(key > node.rgKeys[i]) i++;

```
BTreeInsertNonfull( child, key );
```

4/4

Example BTree:

$$t = 3$$



d.b.s

B-Trees - designed to minimize # times you have to access the disk
internal nodes & ~~external~~/leaf nodes

internal - if it has n keys, it'll have $n+1$ children
relate this to binary trees

Primary vs. Secondary physically moving part

(2) \hookrightarrow cheaper, ^{much} slower; persistent relative to each other
 \hookrightarrow fast, expensive; volatile

HD diagram-

platters on the same spindle

\hookrightarrow covered w/ magnetic material

Arm extends over spindle, head reads \hookrightarrow track

\hookrightarrow all arms are linked

Cylinder - set of tracks - each track divided into a

bunch of pages, each containing a bunch of info.

Disk Access time isn't constant -

Data is too large to put it all into memory

at once - copy little pieces as needed

Alg. mostly bounded by # disk accesses

Min. # disk accesses - make a node as large as a page

B⁺-Tree - store all info (except keys) in leaf nodes - maximizes branching, to minimize # disk accesses

2

Properties - Every node x has - fields -

$N - \#$ keys stored @ x

keys $[]$ - stored in nondecreasing ($\dots \leq \dots$) order

leaf - true if it's a leaf

Internal Nodes -

$x.N + 1$ pointers $x.C[1] \dots x.C[x.N+1]$ to children

The keys separate the ranges of keys stored in each subtree

$$\text{key}_{i-1} \leq k_i \leq \text{key}_i$$

- All leaves have exactly the same depth

\hookrightarrow This is h , the height of the tree

- Upper & lower bounds on # keys a node can contain
bound is $t \geq 2$ min: $\geq t-1$ max: $2t-1$ full: $2t-1$ nodes

2-3-4 trees.

- Every node (except root) must have at least $t-1$ children

\Rightarrow every node (except root) has at least t children

- Every node can contain at most $2t-1$ keys

A node w/ $2t-1$ keys is full

Assume the root stays in-mem - will only need to be written

Find - much like a B.S.T., except that you gotta go through the array of keys till you find the one that's \leq and either that's it, or else, go to child to its immediate left.

(Demo. this briefly)

$$\Theta(h) = \Theta(\log_w n)$$

point of \nearrow Sprawling factor

Create a B-Tree - get the disk page

Mark it as a leaf, w/ 0 entries

Write it to disk

Point root to it.

3

B-Tree Insert

Basic idea: Find a leaf node to insert the value into.

Point out
that you
can store
data (2 nodes,
as well).

B+tree uses
this to maximize
fan-out.

If that node is full (already has $2t - 1$ elts), split that node into 2 separate leaves of $t - 1$ each, with the middle-most being moved up into the parent.

Likewise with the parent, etc.

Once that's done, you can put the new elt. into one node or the other & be sure there's enough space.

Single-Pass Optimization: On your way down the tree, if it's full, split it pre-emptively.

B-Tree-Split-Child(node, i, child)

$z = \text{AllocateNode}()$, // this will be the right $\frac{1}{2}$
copy child's leaf attr to z

$z.\text{numElts} = t - 1$

Copy keys t through $2t - 1$ into z

(if it's not a leaf, copy the child ptrs, too)

$\& \text{child}.\text{numElts} = t - 1$

Shuffle the key/ptr. pairs in node up by one, to make room for z ;

Copy $\text{child}.\text{key}[t]$ into $\text{node}.\text{key}[i]$

$\text{node}.\text{numElts} + 1$

Write node, child, z to disk

This needs a special case for a full root.

4

B-Tree-Insert (Tree, key / data)

if root is full

S = Create a new disk page

S. leaf = false;

Tree.root = S; // save ^{old} root in a temp

S.num Elts = 0

S.Child[1] = r; // BTSC will copy the ptr to the right into S.Child[2]

B-Tree-SplitChild (S, 1, oldRoot)

B-Tree-Insert-Nonfull (Tree, root), key / data)

B-Tree-Insert-Nonfull (node, key / data)

if (node.leaf == true)

{ shuffle all keys up one, put new key into place,
Write node to disk }

}

else // must be @ internal node - figure out where to add

{ search highest to lowest to find its spot

Read the child off the disk.

If (child.num Elts == 2t - 1)

{ B-Tree-Split-Child (node, i, node.child[i])

if (key > x.key[i])

i++;

B-Tree-Insert-Nonfull (node.child[i], key / data)