# B-Trees

**B-trees** are balanced search trees designed to work well on magnetic disks or direct-access secondary storage devices. B-trees differ significantly from red-black trees in that B-tree nodes may have many children, from a handful to thousands.

It often takes more time to access a page of information and read it from a disk than it takes for the computer to examine all the information read. For this reason, we look at the two principal components of the running time:

- the number of disk accesses, and

- the CPU time.

We model these disk operations in our pseudo-code as follows. Let $x$ be a pointer to an object. DISK-READ($x$) to read object $x$ into main memory; DISK-WRITE($x$) to write object $x$ to the disk.

## 1    Definition

A B-tree is a rooted tree having the following properties.

1. Every node $x$ has the following fields:

    - $n[x]$, the number of keys currently stored in node $x$

    - the $n[x]$ keys themselves, stored in nondecreasing order: $key_1[x] \leq key_2[x] \leq \ldots \leq key_{n[x]}[x]$, and
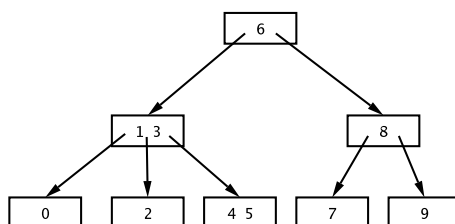


Figure 1: B-Tree where $t = 2$ and the sequence 9,0,8,1,7,2,6,3,5,4 inserted.

- $leaf[x]$, a boolean value that is TRUE if $x$ is a leaf and FALSE if $x$ is an internal node.

2. If $x$ is an internal node, it also contains $n[x]+1$ pointers $c_1[x], c_2[x], \ldots, c_{n[x]+1}[x]$ to its children. Leaf nodes have no children, so their $c_i$ fields are undefined.

3. The keys $key_i[x]$ separate the ranges of keys stored in each subtree: if $k_i$ is any key stored in the subtree with root $c_i[x]$, then

$$k_1 \leq key_1[x] \leq k_2 \leq key_2[x] \leq \ldots \leq key_{n[x]}[x] \leq k_{n[x]+1}.$$

4. Every leaf has the same depth, which is the tree's height $h$.

5. There are lower and upper bounds on the number of keys a node can contain. These bounds can be expressed in terms of a fixed integer $t \geq 2$ called the minimum degree of the B-tree:

   - Every node other than the root must have at least $t-1$ keys. Every internal node other than the root thus has at least $t$ children. If the tree is non-empty, the root must have at least one key.
   - Every node can contain at most $2t-1$ keys. Therefore, an internal node can have at most $2t$ children. We say that a node is **full** if it contains exactly $2t-1$ keys.

**THEOREM** If $n \geq 1$, then for any $n$-key B-tree $T$ of height $h$ and minimum degree $t \geq 2$, then

$$h \leq \log_t \frac{n+1}{2}.$$

# 2  Operations

The operations for B-trees include B-TREE-SEARCH, B-TREE-CREATE, B-TREE-INSERT, B-TREE-DELETE, etc. We assume that:

- the root of the B-tree is always in main memory.

- any nodes that are passed as parameters must already have had a DISK-READ operation performed on them.
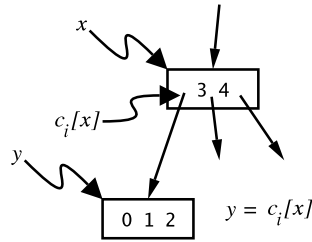
Figure 2: The arguments for the B-Tree-Split-Child algorithm.

## 2.1 Search

B-Tree-Search$(x, k)$
```
1   i ← 1
2   while i ≤ n[x] and k > key_i[x]
3        do i ← i + 1
4   if i ≤ n[x] and k = key_i[x]
5      then return (x, i)
6   if leaf[x]
7      then return NIL
8      else  Disk-Read(c_i[x])
9            return B-Tree-Search(c_i[x], k)
```

## 2.2 Creation

B-Tree-Create$(T)$
```
1   x ← Allocate-Node()
2   leaf[x] ← TRUE
3   n[x] ← 0
4   Disk-Write(x)
5   root[T] ← x
```

## 2.3 Insertion

A fundamental operation used during insertion is the **splitting** of a full node $y$ (having $2t - 1$ keys) around its **median** key $key_t[y]$ into two nodes having $t - 1$ keys each. The median key moves up into $y$'s parent, which must be nonfull prior to the splitting of $y$.
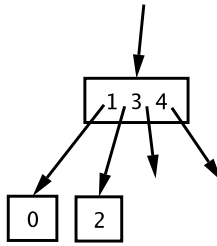
Figure 3: The result of child splitting.

B-Tree-Split-Child$(x, i, y)$
1    $z \leftarrow$ Allocate-Node$()$
2    $leaf[z] \leftarrow leaf[y]$
3    $n[z] \leftarrow t - 1$
4    **for** $j \leftarrow 1$ **to** $t - 1$
5        **do** $key_j[z] \leftarrow key_{j+t}[y]$
6    **if** not $leaf[y]$
7      **then for** $j \leftarrow 1$ **to** $t$
8            **do** $c_j[z] \leftarrow c_{j+t}[y]$
9    $n[y] \leftarrow t - 1$
10   **for** $j \leftarrow n[x] + 1$ **downto** $i + 1$
11      **do** $c_{j+1}[x] \leftarrow c_j[x]$
12   $c_{i+1}[x] \leftarrow z$
13   **for** $j \leftarrow n[x]$ **downto** $i$
14      **do** $key_{j+1}[x] \leftarrow key_j[x]$
15   $key_i[x] \leftarrow key_t[y]$
16   $n[x] \leftarrow n[x] + 1$
17   Disk-Write$(y)$
18   Disk-Write$(z)$
19   Disk-Write$(x)$

Inserting a key $k$ into a B-tree $T$ of height $h$ is done in a single pass down the tree, requiring $O(h)$ disk accesses. The CPU time required is $O(th) = O(t \log_t n)$.

Figure 4: See Fig 18.6 (page 446) – create a new node for root if root is $2t - 1$

B-TREE-INSERT$(T, k)$
1   $r \leftarrow root[T]$
2   **if** $n[r] = 2t - 1$
3       **then** $s \leftarrow$ ALLOCATE-NODE()
4           $root[T] \leftarrow s$
5           $leaf[s] \leftarrow$ FALSE
6           $n[s] \leftarrow 0$
7           $c_1[s] \leftarrow r$
8           B-TREE-SPLIT-CHILD$(s, 1, r)$
9           B-TREE-INSERT-NONFULL$(s, k)$
10      **else** B-TREE-INSERT-NONFULL$(r, k)$

B-TREE-INSERT-NONFULL$(x, k)$
1   $i \leftarrow n[x]$
2   **if** $leaf[x]$
3       **then while** $i \geq 1$ and $k < key_i[x]$
4               **do** $key_{i+1}[x] \leftarrow key_i[x]$
5                   $i \leftarrow i - 1$
6           $key_{i+1}[x] \leftarrow k$
7           $n[x] \leftarrow n[x] + 1$
8           DISK-WRITE$(x)$
9       **else while** $i \geq 1$ and $k < key_i[x]$
10              **do** $i \leftarrow i - 1$
11          $i \leftarrow i + 1$
12          DISK-READ$(c_i[x])$
13          **if** $n[c_i[x]] = 2t - 1$
14              **then** B-TREE-SPLIT-CHILD$(x, i, c_i[x])$
15                  **if** $k > key_i[x]$
16                      **then** $i \leftarrow i + 1$
17          B-TREE-INSERT-NONFULL$(c_i[x], k)$

## 2.4   Deletion

1. If the key $k$ is in node $x$ and $x$ is a leaf, delete the key $k$ from $x$.

2. If the key $k$ is in node $x$ and $x$ is an internal node , do the following:

   - If the child $y$ that precedes $k$ in node $x$ and $x$ has at least $t$ keys, then find the predecessor $k'$ of $k$ in the subtree rooted at $y$. Recursively delete $k'$, and replace $k$ by $k'$ in $x$.

(a) Insert 9, 0, 8

```
0 8 9
```

(b) Insert 1



(c) Insert 7



(d) Insert 2



(e) Insert 6



(f) Insert 3
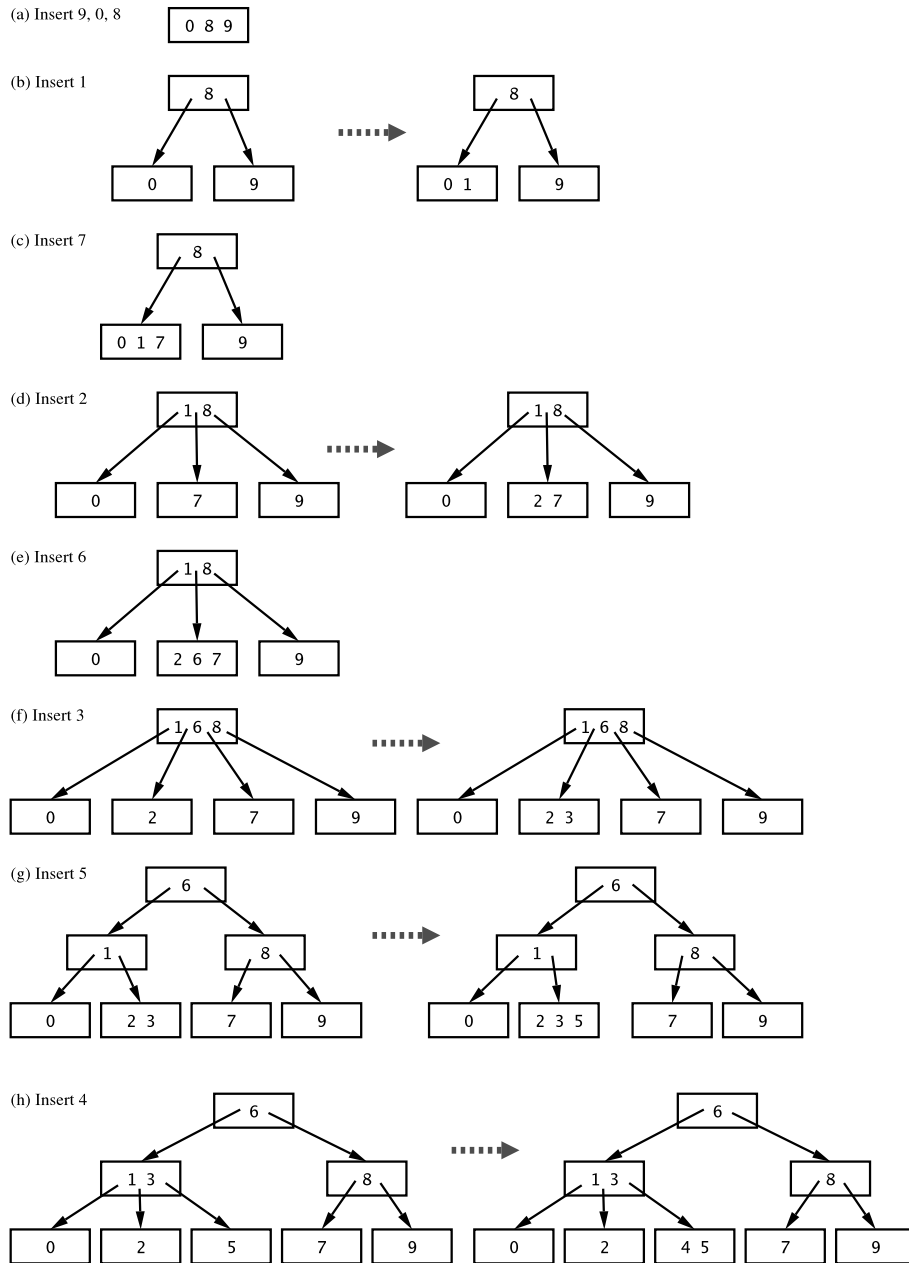


(g) Insert 5



(h) Insert 4



Figure 5: Inserting the sequence 9,0,8,1,7,2,6,3,5,4 into a B-Tree.

- Symmetrically, if the child $z$ that follows $k$ in node $x$ has at least $t$ keys, then find the successor $k'$ of $k$ in the subtree rooted at $z$. Recursively delete $k'$, and replace $k$ by $k'$ in $x$.

- Otherwise, if both $y$ and $z$ have only $t-1$ keys, merge $k$ and all of $z$ into $y$, so that $x$ loses both $k$ and the pointer to $z$, and $y$ now contains $2t-1$ keys. Then free $z$ and recursively delete $k$ from $y$.

3. If the key $k$ is not present in internal node $x$, determine the root $c_i[x]$ of the appropriate subtree that must contain $k$, if $k$ is in the tree at all. If $c_i[x]$ has only $t-1$ keys, execute step 3a or 3b as necessary to guarantee that we descend to a node containing at least $t$ keys. Then, finish by recursing on the appropriate child of $x$.

   (a) If $c_i[x]$ has only $t-1$ keys but has a sibling with $t$ keys, give $c_i[x]$ an extra key by moving a key from $x$ down into $c_i[x]$, moving a key from $c_i[x]$'s immediate left or right sibling up into $x$, and moving the appropriate child from the sibling into $c_i[x]$.

   (b) If $c_i[x]$ and both of $c_i[x]$'s immediate siblings have $t-1$ keys, merge $c_i$ with one sibling., which involves moving a key from $x$ down into the new merged node to become the median key from that node.