

B-Trees

David Kauchak
cs302
Spring 2013



Admin

- Homework 10 out today
- Midterm out Monday/Tuesday
 - Available online
 - 2 hours
 - Will need to return it to me within 3 hours of downloading
 - Must take by Friday at 6pm
- Review on Tuesday
 - E-mail if you have additional topics you'd like covered

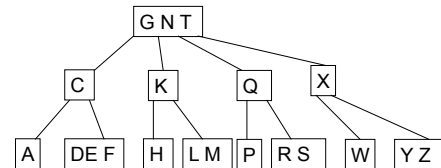


B-tree

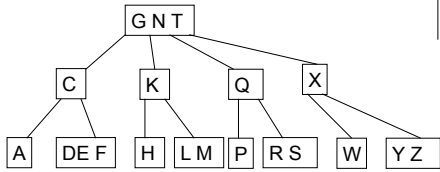
- Defined by one parameter: t
- Balanced n -ary tree
- Each node contains between $t-1$ and $2t-1$ keys/data values (i.e. multiple data values per tree node)
 - keys/data are stored in **sorted order**
 - one exception: root can have $< t-1$ keys
- Each internal node contains between t and $2t$ children
 - the keys of a parent **delimit** the values of the children keys
 - For example, if $key_i = 15$ and $key_{i+1} = 25$ then child $i + 1$ must have keys between 15 and 25
- all leaves have the same depth



Example B-tree: $t = 2$

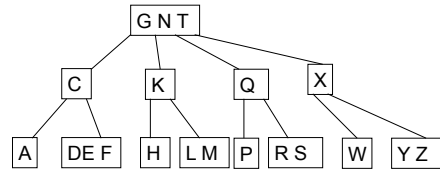


Example B-tree: $t = 2$



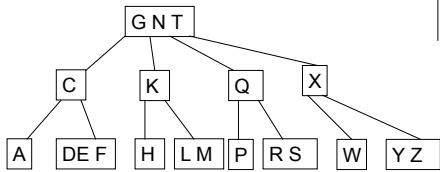
Balanced: all leaves have the same depth

Example B-tree: $t = 2$



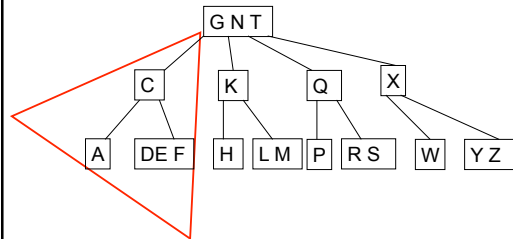
Each node contains between $t-1$ and $2t-1$ keys stored in increasing order

Example B-tree: $t = 2$



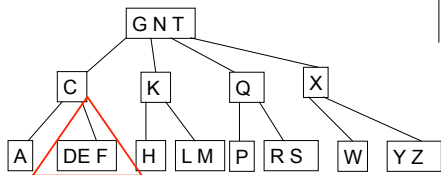
Each node contains between t and $2t$ children

Example B-tree: $t = 2$



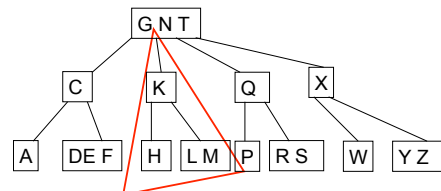
The keys of a parent delimit the values that a child's keys can take

Example B-tree: $t = 2$



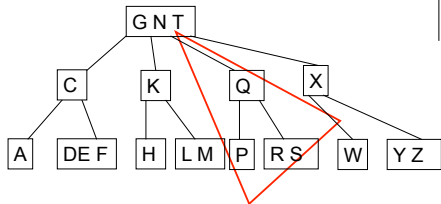
The keys of a parent delimit the values that a child's keys can take

Example B-tree: $t = 2$



The keys of a parent delimit the values that a child's keys can take

Example B-tree: $t = 2$



The keys of a parent delimit the values that a child's keys can take

When do we use B-trees over other balanced trees?

B-trees are generally an **on-disk** data structure

Memory is limited or there is a large amount of data to be stored

In the extreme, only one node is kept in memory and the rest on disk

Size of the nodes is often determined by a page size on disk. **Why?**

Databases frequently use B-trees

Notes about B-trees

Because t is generally large, the height of a B-tree is usually small

$t = 1001$ with height 2, how many values can we have?

Each internal node contains between t and $2t$ children

Each node contains between $t-1$ and $2t-1$ keys/data values (i.e. multiple data values per tree node)

root level 1 level 2

$$2001 + 2002 * 2001 + 2002 * 2002 * 2001 = 8,024,024,007$$

(over 8 billion keys!!!)

Notes about B-trees

Because t is generally large, the height of a B-tree is usually small

We will count both run-time as well as the number of disk accesses. **Why?**

Height of a B-tree

B-trees have a similar feeling to BSTs

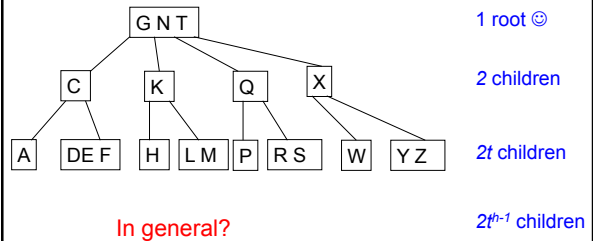
We saw for BSTs that most of the operations depended on the height of the tree

How can we bound the height of the tree?

We know that nodes must have a minimum number of keys/data items ($t-1$)

For a tree of height h , what is the smallest number of keys?

Minimum number of nodes at each depth?



Minimum number of keys/values

$$n \geq 1 + (t-1) \sum_{i=1}^h 2t^{i-1}$$

↑ root ↑ min. keys per node ↑ min. number of nodes

Minimum number of nodes

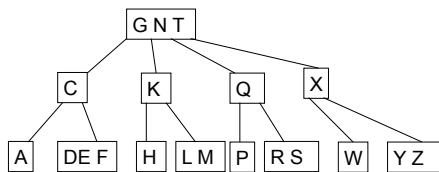
$$\begin{aligned}
 n &\geq 1 + (t-1) \sum_{i=1}^h 2t^{i-1} \\
 &= 1 + 2(t-1) \left(\frac{t^h - 1}{t-1} \right) \\
 &= 2t^h - 1
 \end{aligned}$$

so,

$$\begin{aligned}
 t^h &\leq (n+1)/2 \\
 h &\leq \log_t \frac{(n+1)}{2}
 \end{aligned}$$

Searching B-Trees

Find value k in B-Tree



Searching B-Trees

Find value k in B-Tree node x

```

B-TREESEARCH(x, k)
1  i ← 1
2  while i ≤ n(x) and k > Kx[i]
3    i ← i + 1
4  if i ≤ n(x) and k = Kx[i]
5    return (x, i)
6  if LEAF(x)
7    return null
8  else
9    DISKREAD(Cx[i])
10   return B-TREESEARCH(Cx[i], k)
    
```

↖ number of keys ↖ key[i] ↖ child[i]

Searching B-Trees

B-TREESEARCH(x, k)

```

1   $i \leftarrow 1$ 
2  while  $i \leq n(x)$  and  $k > K_x[i]$ 
3       $i \leftarrow i + 1$ 
4  if  $i \leq n(x)$  and  $k = K_x[i]$ 
5      return  $(x, i)$ 
6  if LEAF( $x$ )
7      return null
8  else
9      DISKREAD( $C_x[i]$ )
10     return B-TREESEARCH( $C_x[i], k$ )

```

make disk reads
explicit

Searching B-Trees

B-TREESEARCH(x, k)

```

1   $i \leftarrow 1$ 
2  while  $i \leq n(x)$  and  $k > K_x[i]$ 
3       $i \leftarrow i + 1$ 
4  if  $i \leq n(x)$  and  $k = K_x[i]$ 
5      return  $(x, i)$ 
6  if LEAF( $x$ )
7      return null
8  else
9      DISKREAD( $C_x[i]$ )
10     return B-TREESEARCH( $C_x[i], k$ )

```

iterate through the sorted keys
and find the correct location

Searching B-Trees

B-TREESEARCH(x, k)

```

1   $i \leftarrow 1$ 
2  while  $i \leq n(x)$  and  $k > K_x[i]$ 
3       $i \leftarrow i + 1$ 
4  if  $i \leq n(x)$  and  $k = K_x[i]$ 
5      return  $(x, i)$ 
6  if LEAF( $x$ )
7      return null
8  else
9      DISKREAD( $C_x[i]$ )
10     return B-TREESEARCH( $C_x[i], k$ )

```

if we find the value
in this node, return it

Searching B-Trees

B-TREESEARCH(x, k)

```

1   $i \leftarrow 1$ 
2  while  $i \leq n(x)$  and  $k > K_x[i]$ 
3       $i \leftarrow i + 1$ 
4  if  $i \leq n(x)$  and  $k = K_x[i]$ 
5      return  $(x, i)$ 
6  if LEAF( $x$ )
7      return null
8  else
9      DISKREAD( $C_x[i]$ )
10     return B-TREESEARCH( $C_x[i], k$ )

```

if it's a leaf and we didn't
find it, it's not in the tree

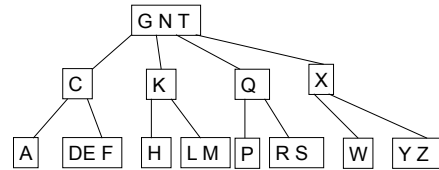
Searching B-Trees

```

B-TREESEARCH(x, k)
1  i ← 1
2  while i ≤ n(x) and k > Kx[i]
3      i ← i + 1
4  if i ≤ n(x) and k = Kx[i]
5      return (x, i)
6  if LEAF(x)
7      return null
8  else
9      DISKREAD(Cx[i])
10     return B-TREESEARCH(Cx[i], k)
    
```

Recurse on the proper child where the value is between the keys

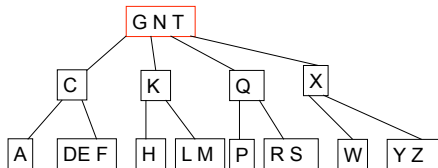
Search example: R



```

B-TREESEARCH(x, k)
1  i ← 1
2  while i ≤ n(x) and k > Kx[i]
3      i ← i + 1
4  if i ≤ n(x) and k = Kx[i]
5      return (x, i)
6  if LEAF(x)
7      return null
8  else
9      DISKREAD(Cx[i])
10     return B-TREESEARCH(Cx[i], k)
    
```

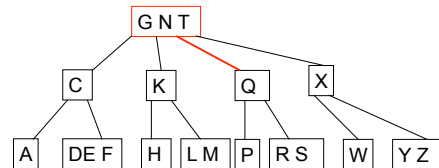
Search example: R



```

B-TREESEARCH(x, k)
1  i ← 1
2  while i ≤ n(x) and k > Kx[i]
3      i ← i + 1
4  if i ≤ n(x) and k = Kx[i]
5      return (x, i)
6  if LEAF(x)
7      return null
8  else
9      DISKREAD(Cx[i])
10     return B-TREESEARCH(Cx[i], k)
    
```

Search example: R



```

B-TREESEARCH(x, k)
1  i ← 1
2  while i ≤ n(x) and k > Kx[i]
3      i ← i + 1
4  if i ≤ n(x) and k = Kx[i]
5      return (x, i)
6  if LEAF(x)
7      return null
8  else
9      DISKREAD(Cx[i])
10     return B-TREESEARCH(Cx[i], k)
    
```

find the correct location

Search example: R

```

B-TREESEARCH(x, k)
1 i ← 1
2 while i ≤ n(x) and k > Kx[i]
3   i ← i + 1
4 if i ≤ n(x) and k = Kx[i]
5   return (x, i)
6 if LEAF(x)
7   return null
8 else
9   DISKREAD(Cx[i])
10  return B-TREESEARCH(Cx[i], k)
    
```

the value is not in this node

Search example: R

```

B-TREESEARCH(x, k)
1 i ← 1
2 while i ≤ n(x) and k > Kx[i]
3   i ← i + 1
4 if i ≤ n(x) and k = Kx[i]
5   return (x, i)
6 if LEAF(x)
7   return null
8 else
9   DISKREAD(Cx[i])
10  return B-TREESEARCH(Cx[i], k)
    
```

this is not a leaf node

Search example: R

```

B-TREESEARCH(x, k)
1 i ← 1
2 while i ≤ n(x) and k > Kx[i]
3   i ← i + 1
4 if i ≤ n(x) and k = Kx[i]
5   return (x, i)
6 if LEAF(x)
7   return null
8 else
9   DISKREAD(Cx[i])
10  return B-TREESEARCH(Cx[i], k)
    
```

Search example: R

```

B-TREESEARCH(x, k)
1 i ← 1
2 while i ≤ n(x) and k > Kx[i]
3   i ← i + 1
4 if i ≤ n(x) and k = Kx[i]
5   return (x, i)
6 if LEAF(x)
7   return null
8 else
9   DISKREAD(Cx[i])
10  return B-TREESEARCH(Cx[i], k)
    
```

find the correct location

Search example: R

```

B-TREESEARCH(x, k)
1  i ← 1
2  while i ≤ n(x) and k > Kx[i]
3    i ← i + 1
4  if i ≤ n(x) and k = Kx[i]
5    return (x, i)
6  if LEAF(x)
7    return null
8  else
9    DISKREAD(Cx[i])
10   return B-TREESEARCH(Cx[i], k)
    
```

not in this node and this is not a leaf

Search example: R

```

B-TREESEARCH(x, k)
1  i ← 1
2  while i ≤ n(x) and k > Kx[i]
3    i ← i + 1
4  if i ≤ n(x) and k = Kx[i]
5    return (x, i)
6  if LEAF(x)
7    return null
8  else
9    DISKREAD(Cx[i])
10   return B-TREESEARCH(Cx[i], k)
    
```

Search example: R

```

B-TREESEARCH(x, k)
1  i ← 1
2  while i ≤ n(x) and k > Kx[i]
3    i ← i + 1
4  if i ≤ n(x) and k = Kx[i]
5    return (x, i)
6  if LEAF(x)
7    return null
8  else
9    DISKREAD(Cx[i])
10   return B-TREESEARCH(Cx[i], k)
    
```

find the correct location

Search example: R

```

B-TREESEARCH(x, k)
1  i ← 1
2  while i ≤ n(x) and k > Kx[i]
3    i ← i + 1
4  if i ≤ n(x) and k = Kx[i]
5    return (x, i)
6  if LEAF(x)
7    return null
8  else
9    DISKREAD(Cx[i])
10   return B-TREESEARCH(Cx[i], k)
    
```

Search running time

How many calls to BTreeSearch?

- $O(\text{height of the tree})$
- $O(\log n)$

Disk accesses?

- One for each call – $O(\log n)$

Computational time?

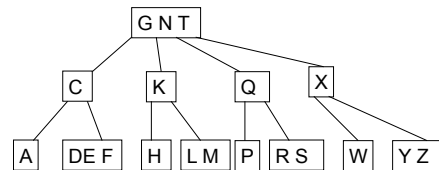
- $O(t)$ keys per node
- linear search
- $O(t \log n)$

Why not binary search to find key in a node?

```

B-TREESEARCH(x, k)
1  i ← 1
2  while i ≤ n(x) and k > Ki[t]
3    i ← i + 1
4  if i ≤ n(x) and k = Ki[t]
5    return (x, i)
6  if LEAF(x)
7    return null
8  else
9    DISKREAD(Ci[t])
10   return B-TREESEARCH(Ci[t], k)
    
```

BST-Insert



B-Tree insert

Starting at root, follow the *search* path down the tree

- If the node is full (contains $2t - 1$ keys)
 - split the keys into two nodes around the median value
 - add the median value to the parent node
- If the node is a leaf, insert it into the correct spot

Observations

- Insertions **always** happens in the leaves
- When does the height of a B-tree grow?
- Why do we know it's always ok when we're splitting a node to insert the median value into the parent?

Insertion: $t = 2$

G C N A H E K Q M F W L T Z D P R X Y S

Insertion: t = 2

GCNAHEKQMF~~W~~LZDPRXYS

G



Insertion: t = 2

GCNAHEKQMF~~W~~LZDPRXYS


CG



Insertion: t = 2

GCNAHEKQMF~~W~~LZDPRXYS

CGN




Insertion: t = 2

GCNAHEKQMF~~W~~LZDPRXYS

CGN

Node is full, so split



Insertion: t = 2

GCNAHEKQMFWLTZDPRXYS

Node is full, so split

Insertion: t = 2

GCNAHEKQMFWLTZDPRXYS

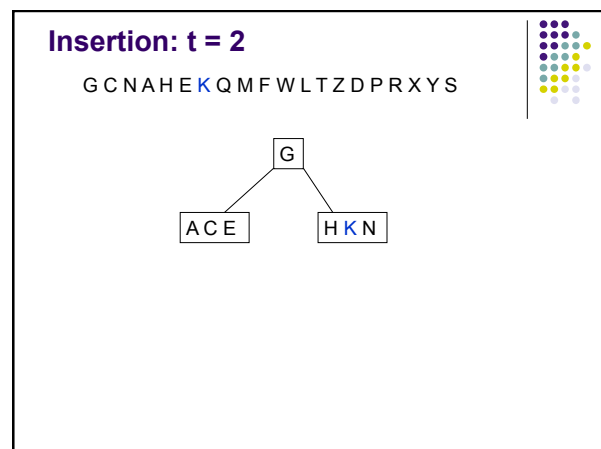
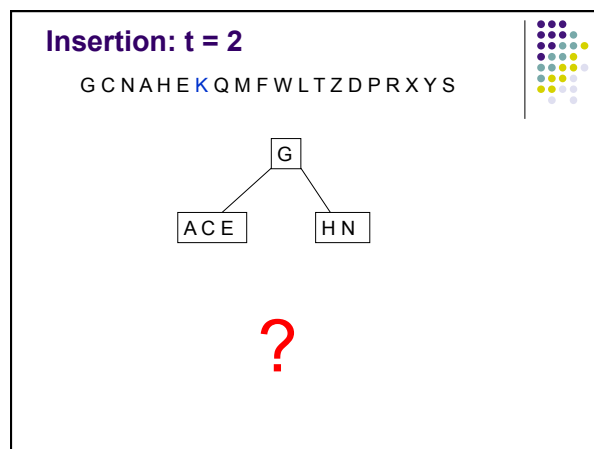
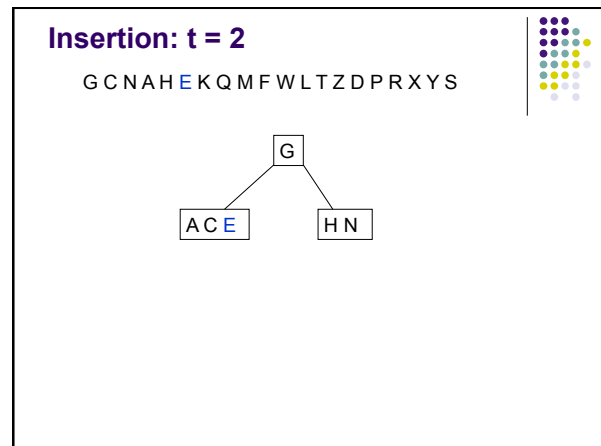
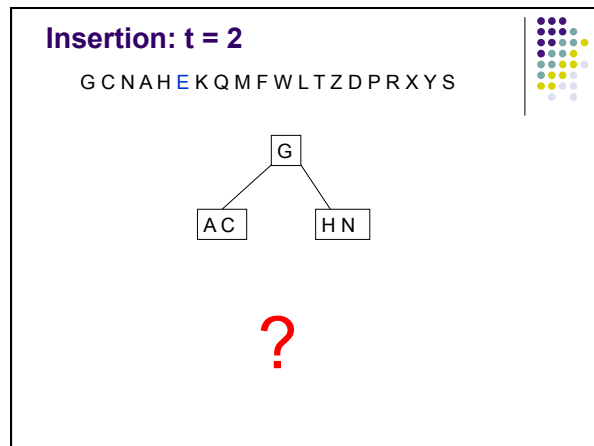
Insertion: t = 2

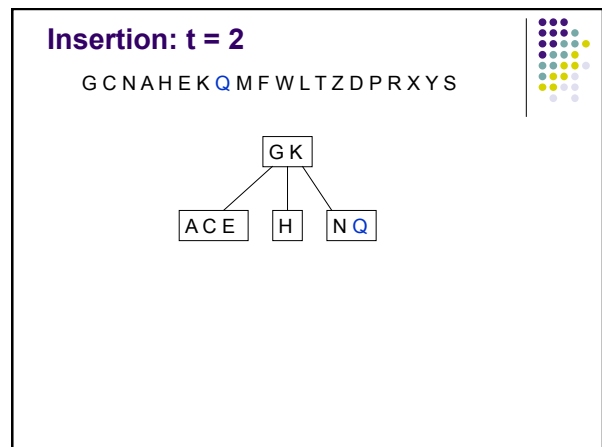
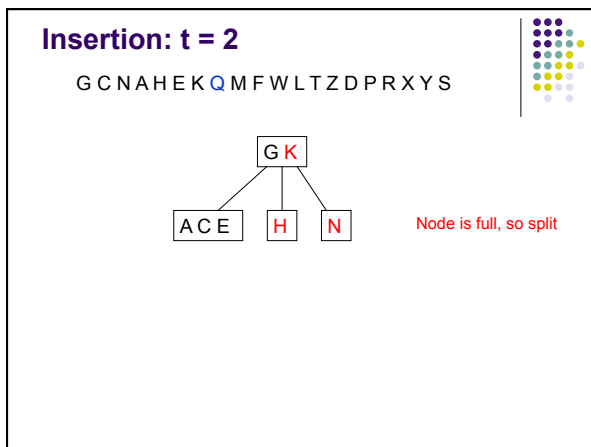
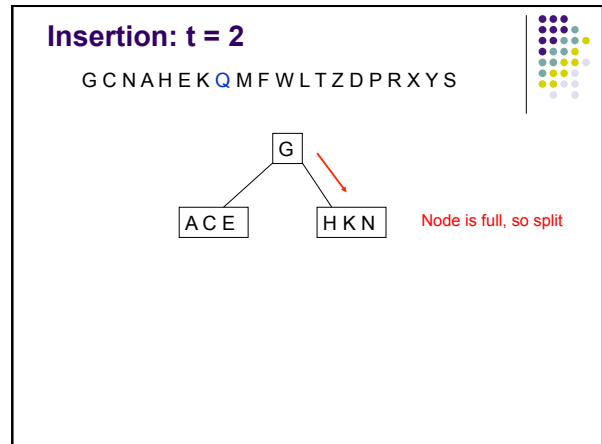
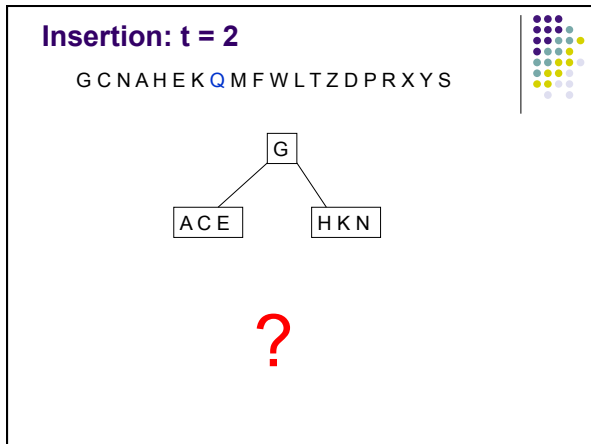
GCNAHEKQMFWLTZDPRXYS

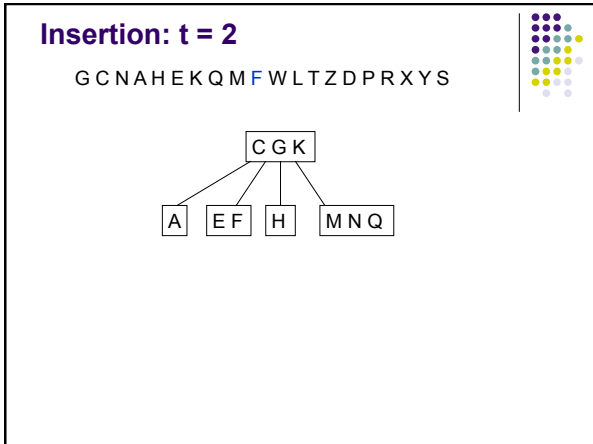
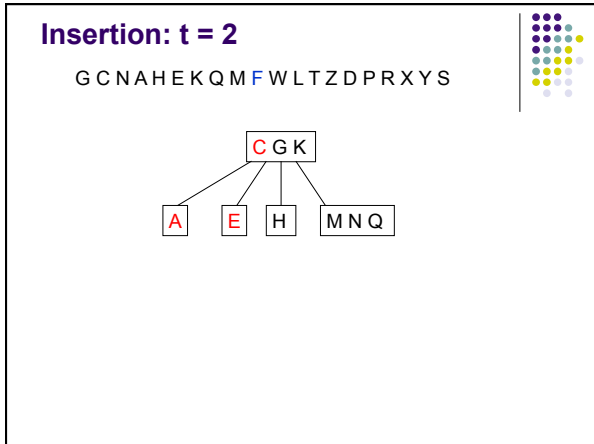
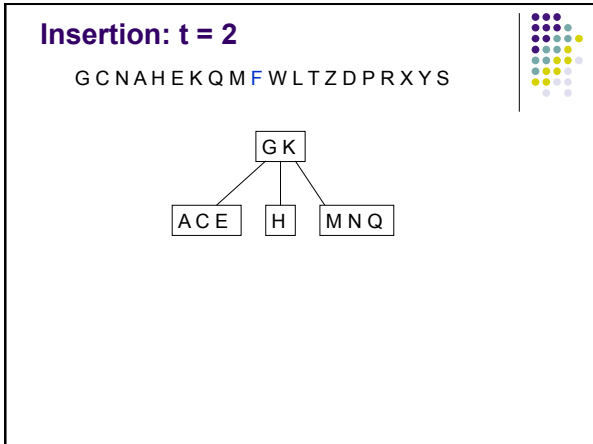
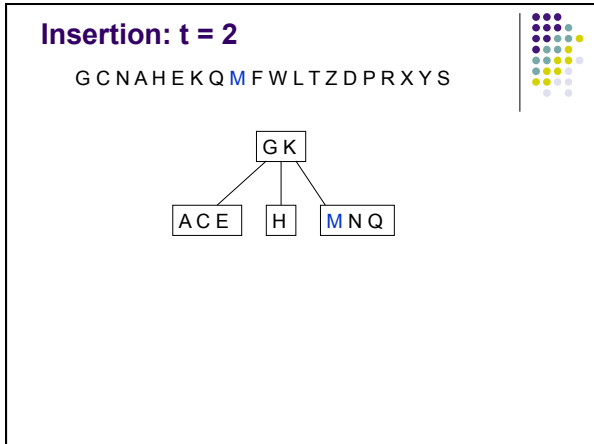
?

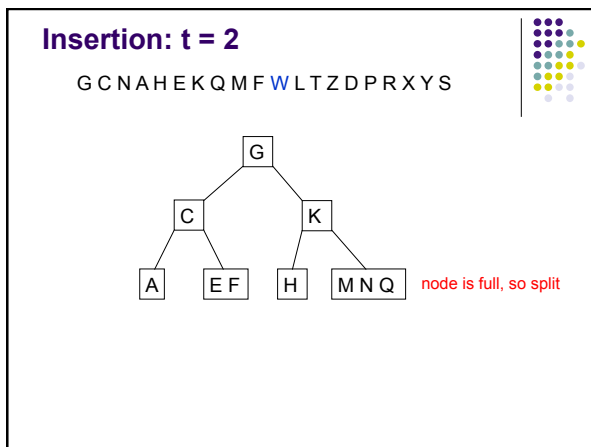
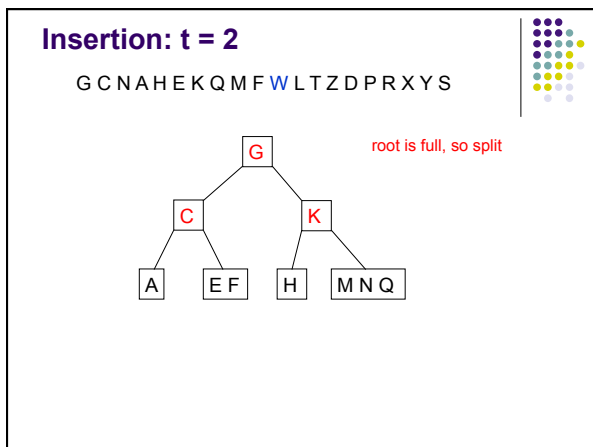
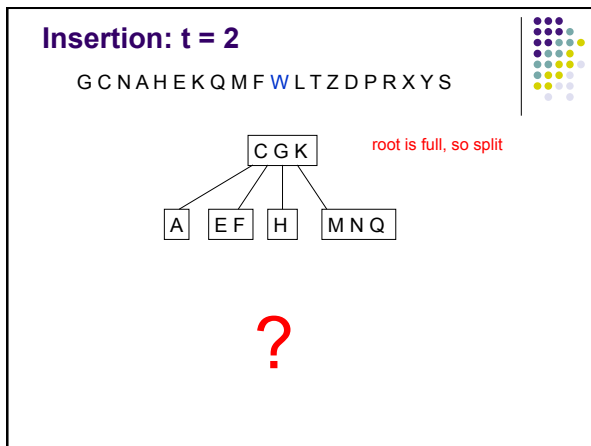
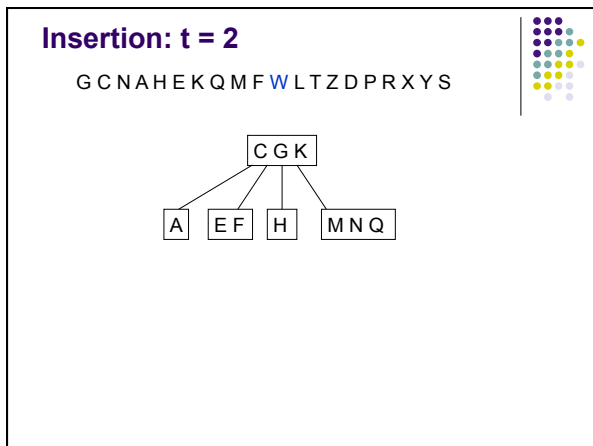
Insertion: t = 2

GCNAHEKQMFWLTZDPRXYS









Insertion: $t = 2$
 GCNAHEKQMF~~W~~LTZDPRXYS

node is full, so split

Insertion: $t = 2$
 GCNAHEKQMF~~W~~LTZDPRXYS

Insertion: $t = 2$
 GCNAHEKQMF~~W~~...

Correctness of insert

Starting at root, follow *search* path down the tree

- If the node is full (contains $2t - 1$ keys), split the keys around the median value into two nodes and add the median value to the parent node
- If the node is a leaf, insert it into the correct spot

Does it add the value in the correct spot?

- Follows the correct *search* path
- Inserts in correct position

Correctness of insert



Starting at root, follow *search* path down the tree

- If the node is full (contains $2t - 1$ keys), split the keys around the median value into two nodes and add the median value to the parent node
- If the node is a leaf, insert it into the correct spot

Do we maintain a proper B-tree?

- **Maintain $t-1$ to $2t-1$ keys per node?**
 - Always split full nodes when we see them
 - Only split full nodes
- **All leaves at the same level?**
 - Only add nodes at leaves

Insert running time



Without any splitting?

- Similar to BTreeSearch, with one extra disk write at the leaf
- $O(\log_r n)$ disk accesses
- $O(t \log_r n)$ computation time

When a node is split



How many disk accesses?

- 3 disk write operations
 - 2 for the new nodes created by the split (one is reused, but must be updated)
 - 1 for the parent node to add median value

Runtime to split a node?

- $O(t)$ – iterating through the elements a few times since they're already in sorted order

Maximum number of nodes split for a call to insert?

- $O(\text{height of the tree})$

Running time of insert



$O(\log_r n)$ disk accesses

$O(t \log_r n)$ computational costs

Deleting a node from a B-tree



Similar to insertion

- must make sure we maintain B-tree properties (i.e. all leaves same depth and key/node restrictions)
- Proactively move a key from a child to a parent if the parent has $t-1$ keys

$O(\log, n)$ disk accesses

$O(t \log, n)$ computational costs

Summary of operations



Search, Insertion, Deletion

- disk accesses: $O(\log, n)$
- computation: $O(t \log, n)$

Max, Min

- disk accesses: $O(\log, n)$
- computation: $O(\log, n)$

Tree traversal

- disk accesses: if $2t \sim$ page size: $O(\text{minimum \# pages to store data})$
- Computation: $O(n)$