

Algorithms

Chapter 18

B-Trees

- ①可以控制高度的搜尋樹
- ②配合磁碟的特性做最佳化

Associate Professor: Ching-Chi Lin

林清池 副教授

chingchi.lin@gmail.com

Department of Computer Science and Engineering
National Taiwan Ocean University

Outline

- ▶ **Definition of B-trees** B-tree 的性质
- ▶ Basic operations on B-trees search + insert
- ▶ Deleting a key from a B-tree delete

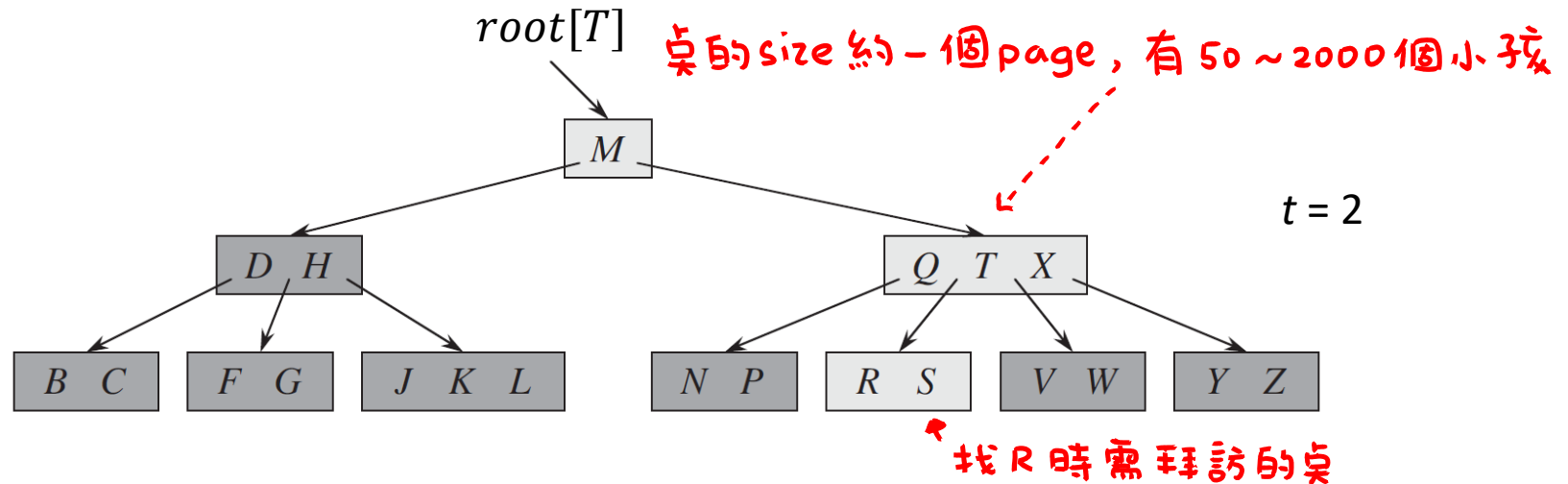
Overview_{1/2}

紅黑樹: 高度 $O(\lg n)$
insert, delete, search 時間 $O(\lg n)$

B-trees: 高度 $O(\log_t n)$ 孩子的個數 $t \sim 2t$
insert, delete, search 的時間 $O(t \log_t n)$

- ▶ B-trees are balanced search trees designed to work well on magnetic disks. 配合磁碟的特性做最佳化
- ▶ B-trees are similar to red-black trees, but they are better at minimizing disk I/O operations.
- ▶ **Red-black trees** = searchs + balanced
 - ▶ A variation of binary search trees.
 - ▶ **Balanced**: height is $O(\lg n)$, where n is the number of nodes.
 - ▶ Operations will take $O(\lg n)$ time in the worst case.
- ▶ **B-trees** = search tree + balanced + magnetic disks
 - ▶ Generalize binary search trees in a natural manner.
 - ▶ **Balanced**: height is $O(\lg_t n)$, where n is the number of nodes.
 - ▶ Operations will take $O(t \lg_t n)$ time in the worst case, where t is the minimum degree of the B-tree.

Overview_{2/2}



- ▶ Running time of a B-tree algorithm is determined by the number of DISK-READ and DISK-WRITE operations.
 - ▶ Thus, a B-tree node is usually as large as a whole disk page.
 - ▶ The "branching factors" between 50 and 2000 are often used, depending on the size of a key relative to the size of a page.
- ▶ An internal node x containing $n[x]$ keys has $n[x]+1$ children.

有 $n[x]$ 個 keys \Rightarrow 有 $n[x]+1$ 個兒子

重點1 \Rightarrow 非 root, key 的個數要在 $t-1 \leq n[x] \leq 2t-1$ 之間
是 root, key 的個數要在 $1 \leq n[x] \leq 2t-1$ 之間

Properties of B-trees_{1/2}

► A **B-tree** T is a rooted tree having the following properties:

► Every node x has the following fields:

► $n[x]$, the number of keys in node x ,

► $key_1[x] \leq key_2[x] \leq \dots \leq key_{n[x]}[x]$,

► $leaf[x]$, $leaf[x] = \text{TRUE}$ if x is a leaf, $leaf[x] = \text{FALSE}$ if x is an internal node.

$n[x]$ = key 值的個數

小 $\xrightarrow{\hspace{1cm}}$ 大

► Each internal node x also contains $n[x]+1$ pointers $c_1[x], c_2[x], \dots, c_{n[x]+1}[x]$ to its children.

► If k_i is any key stored in the subtree with root $c_i[x]$, then

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

重點2 \Rightarrow 所有 leaves 有
相同深度

► All leaves have the same depth, which is the tree's height h .

► Every node x other than the root must have $t-1 \leq n[x] \leq 2t-1$, where $t \geq 2$ is the **minimum degree** of the B-tree.

► If the tree is nonempty, the root has $1 \leq n[root] \leq 2t-1$.

Properties of B-trees_{2/2}

- ▶ The simplest B-tree occurs when $t = 2$. 最簡單的 B-tree, 当 $t = 2$ 時
- ▶ Every internal node then has either 2, 3, or 4 children, and we have a **2-3-4 tree**. $t = 2$ 時, 它可能有 2~4 個兒子
⇒ 2, 3, 4 樹

樹高: $O(\log_t n)$

Height of a B-tree $\log_t n = \frac{\lg n}{\lg t}$, 樹高較紅黑樹小 $\lg t$ 倍

- **Lemma 18.1** If $n \geq 1$, then for any n -key B-tree T of height h and minimum degree $t \geq 2$, $h \leq \log_t \frac{n+1}{2}$.

Proof:

- The root contains at least one key.
- Thus, there are at least 2 nodes at depth 1.
- All other nodes contain at least $t - 1$ keys.
- So, at least $2t$ nodes at depth 2, at least $2t^2$ nodes at depth 3, and so on.

► Then, we have $n \geq \overset{\text{第0層}}{\underset{\downarrow}{1}} + (t-1) \sum_{i=1}^h 2t^{i-1} \leftarrow \begin{matrix} \text{第1} \sim \text{第} h \text{層} \\ \text{每個父至少有 } t-1 \text{ 個鍵值} \end{matrix}$

$$= 1 + 2(t-1) \left(\frac{t^h - 1}{t - 1} \right)$$
$$= 2t^h - 1.$$

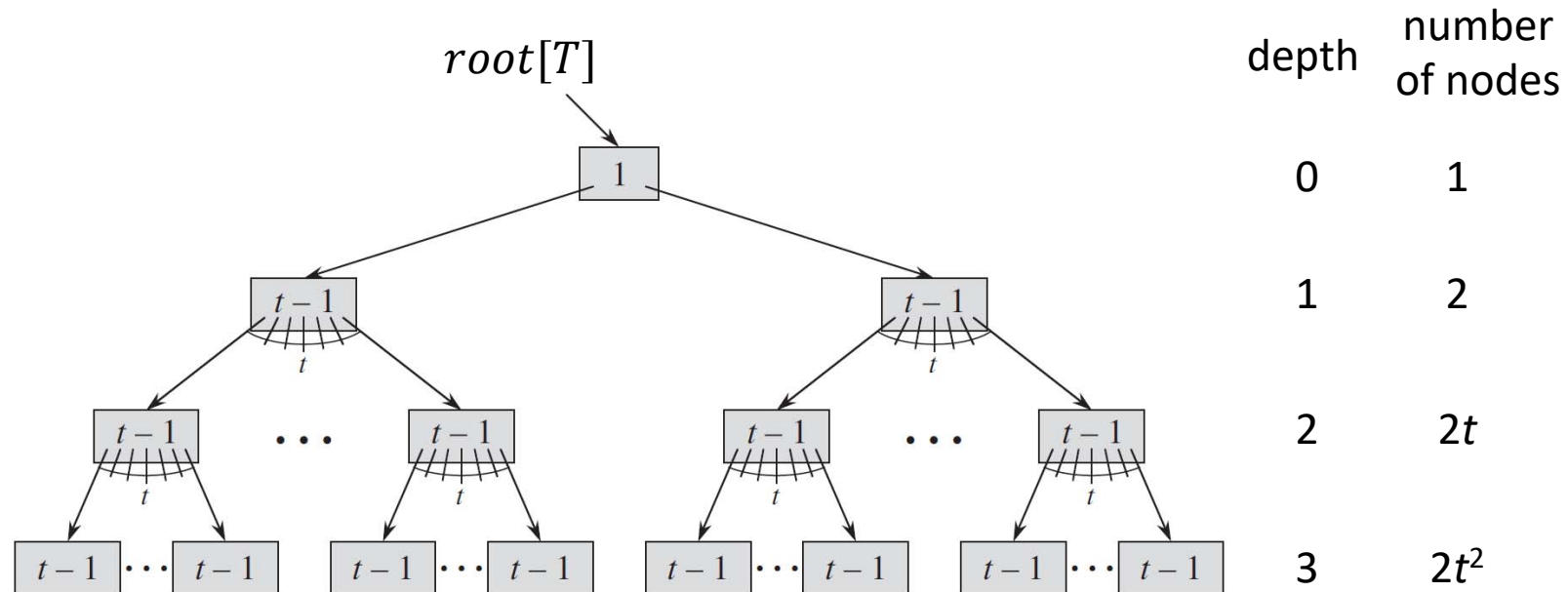
$h \leq \log_t \frac{n+1}{2}$

想法：樹要高

⇒ 每一層的葉數愈少愈好

⇒ 每一個葉的 key 值愈少愈好

Height of a B-tree_{2/2}



第0層：是 root，只有一葉，最少有一個 key 值

第1層：因為 root 只有一個 key 值至少有 2 葉，

第2層：第1層每一個葉至少有 $t-1$ 個 key 值，
也就是 t 個兒子，故至少有 $2t$ 個葉

Outline

- ▶ Definition of B-trees
- ▶ **Basic operations on B-trees** 搜尋與新增
- ▶ Deleting a key from a B-tree

Searching a B-tree 与 binary tree search 相似, 只是要比 $n[x]$ 個 鍵值

- ▶ **B-TREE-SEARCH** is a straightforward generalization of the **TREE-SEARCH** procedure defined for binary search trees.
- ▶ Instead of making a binary, or “two-way” branching decision at each node, we make an $(n[x]+1)$ -way branching decision.

B-TREE-SEARCH(x, k)

```
1.   $i \leftarrow 1$ 
2.  while  $i \leq n[x]$  and  $k > key_i[x]$ 
3.       $i \leftarrow i + 1$ 
4.  if  $i \leq n[x]$  and  $k = key_i[x]$ 
5.      return  $(x, i)$ 
6.  elseif  $leaf[x]$ 
7.      return NIL
8.  else DISK-READ( $c_i[x]$ )
9.      return B-TREE-SEARCH( $c_i[x], k$ )
```

小 大
由左到右 - 一個一個比
相同: 找到了
leaf: 不存在
往子樹, 繼續找
時間: $O(t \log_t n)$

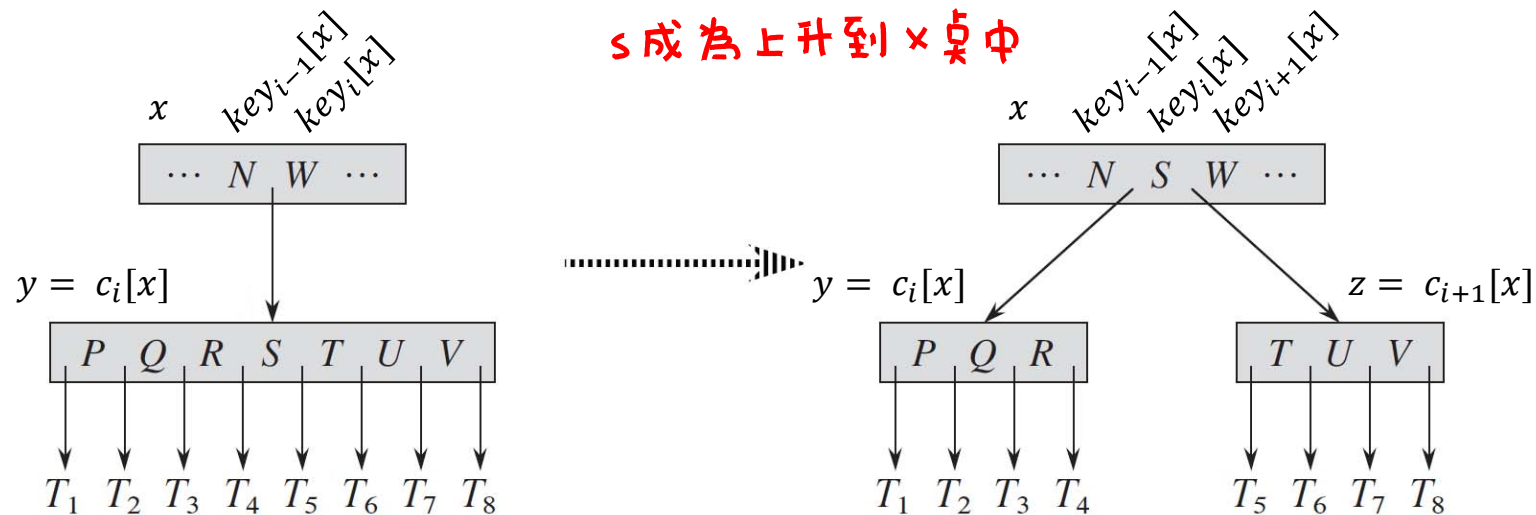
- ▶ Since $n[x] < 2t$, the running is $O(th) = O(t \log_t n)$.

insert 的过程中可能需做 splitting 的动作

將一個有 $2t-1$ 個鍵值的桌分割成 2 個有 $t-1$ 個鍵值的桌

Splitting a node in a B-tree_{1/2}

- ▶ A fundamental operation used during insertion is the **splitting** a full node y (having $2t - 1$ keys) around its **median key** $key_t[y]$ into two nodes having $t - 1$ keys each.



Splitting a node with $t = 4$. Node y is split into two nodes, y and z , and the median key S of y is moved up into y 's parent.

If y has no parent, then the tree grows in height by one.

Splitting a node in a B-tree_{2/2}

B-TREE-SPLIT-CHILD(x, i, y)

y 是 x 的第 i 個兒子

```
1.   $z \leftarrow \text{ALLOCATE-NODE}()$ 
2.   $\text{leaf}[z] \leftarrow \text{leaf}[y]$ 
3.   $n[z] \leftarrow t - 1$ 
4.  for  $j \leftarrow 1$  to  $t - 1$ 
5.       $\text{key}_j[z] \leftarrow \text{key}_{j+t}[y]$ 
6.  if not  $\text{leaf}[y]$ 
7.      for  $j \leftarrow 1$  to  $t$ 
8.           $c_j[z] \leftarrow c_{j+t}[y]$ 
9.   $n[y] \leftarrow t - 1$ 
10. for  $j \leftarrow n[x] + 1$  downto  $i + 1$ 
11.      $c_{j+1}[x] \leftarrow c_j[x]$ 
12.  $c_{i+1}[x] \leftarrow z$ 
13. for  $j \leftarrow n[x]$  downto  $i$ 
14.      $\text{key}_{j+1}[x] \leftarrow \text{key}_j[x]$ 
15.  $\text{key}_i[x] \leftarrow \text{key}_t[y]$ 
16.  $n[x] \leftarrow n[x] + 1$ 
17. DISK-WRITE( $y$ )
18. DISK-WRITE( $z$ )
19. DISK-WRITE( $x$ )
```

$\Theta(1)$ 1~3 建立 z 節

4~8 建立 z 的鍵值和兒子

$\Theta(t)$ 9 設定 y 有 $t-1$ 個兒子

10~14 因為 y 一分為二作調整

① 將 x 的 key 和 child 指標右移

② 放入 x 的第 i 個兒子

③ 設定 x 有 $n[x]+1$ 個兒子

$\Theta(1)$

Time : $O(t)$

Binary search tree: 從 root 往下找一個合適的父親

B-tree: 作法相同, 但保證尋找過程中的節都還沒滿

Inserting a key into a B-tree_{1/3}

- ▶ **B-TREE-INSERT** inserts a key k into a B-tree T of height h in a single pass down the tree. Requiring $O(h)$ disk accesses.
- ▶ Use B-TREE-SPLIT-CHILD to guarantee that the recursion never descends to a full node. 使用 splitting 保證往下的過程中都不會滿

B-TREE-INSERT(T, k)

1. $r \leftarrow \text{root}[T]$
2. if $n[r] = 2t - 1$
3. $s \leftarrow \text{ALLOCATE-NODE}()$
4. $\text{root}[T] \leftarrow s$
5. $\text{leaf}[s] \leftarrow \text{FALSE}$
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)

先檢查 root 有沒有滿
如果有, 將 root 分為二

Handle the case in which the root node r is full: the root is split and a new node s becomes the root.

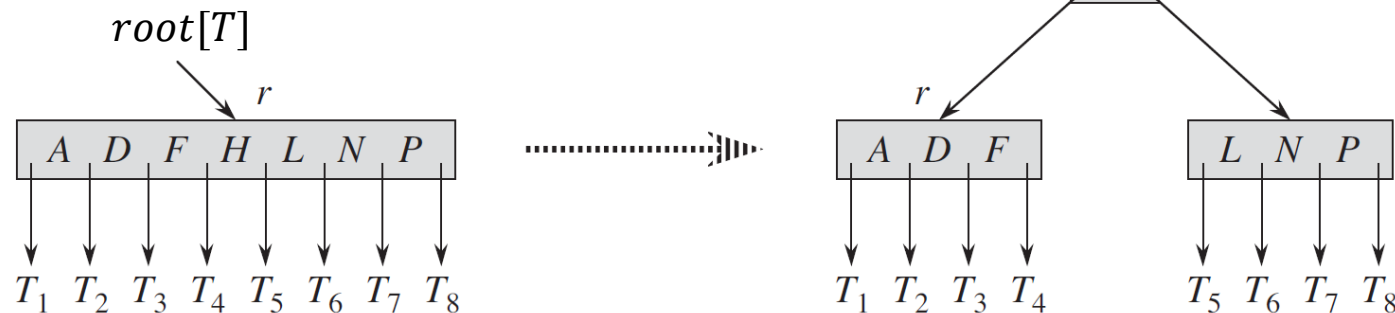
已經確認 root 沒滿

- ▶ The CPU time required is $O(th) = O(t \log_t n)$. 時間: $O(t \log_t n)$

Inserting a key into a B-tree_{2/3}

- ▶ Unlike a binary search tree, a B-tree increases in height at the top instead of at the bottom.

B-tree 在 top 長高
Binary search tree 在 bottom 長高



Splitting the root with $t = 4$. Root node r is split in two, and a new root node s is created.

The B-tree grows in height by one when the root is split.

如果 root 分割，則 B-tree 的高度加 1

Inserting a key into a B-tree_{3/3}

- ▶ B-TREE-INSERT-NONFULL inserts key k into node x , which is assumed to be nonfull when the procedure is called.

B-TREE-INSERT-NONFULL(x, k) 已知 x 沒有滿

```
1.   $i \leftarrow n[x]$ 
2.  if leaf[ $x$ ]
3.      while  $i \geq 1$  and  $k < key_i[x]$ 
4.           $key_{i+1}[x] \leftarrow key_i[x]$ 
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.       $i \leftarrow i - 1$ 
11.       $i \leftarrow i + 1$ 
12.      DISK-READ( $c_i[x]$ )
13.      if  $n[c_i[x]] = 2t - 1$ 
14.          B-TREE-SPLIT-CHILD( $x, i, c_i[x]$ )
15.          if  $k > key_i[x]$ 
16.               $i \leftarrow i + 1$ 
17.          B-TREE-INSERT-NONFULL( $c_i[x], k$ )
```

小

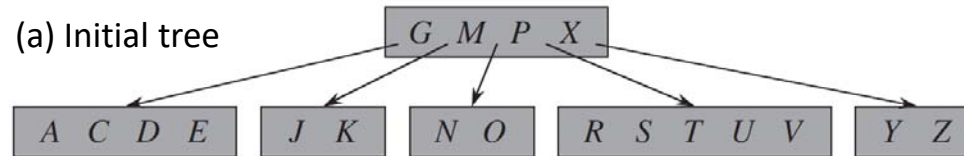
大



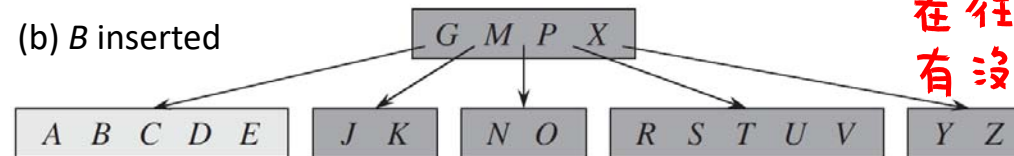
leaf: 從右往左比, 找合適的地方放

internal: 從右往左比, 找合適的子樹,
在往下之前先檢查兒子
有沒有滿

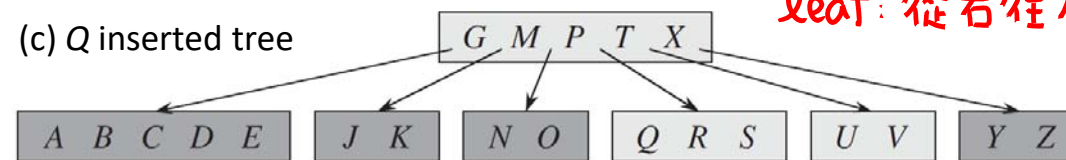
Time : $O(th)$
 $= O(t \log_t n)$.



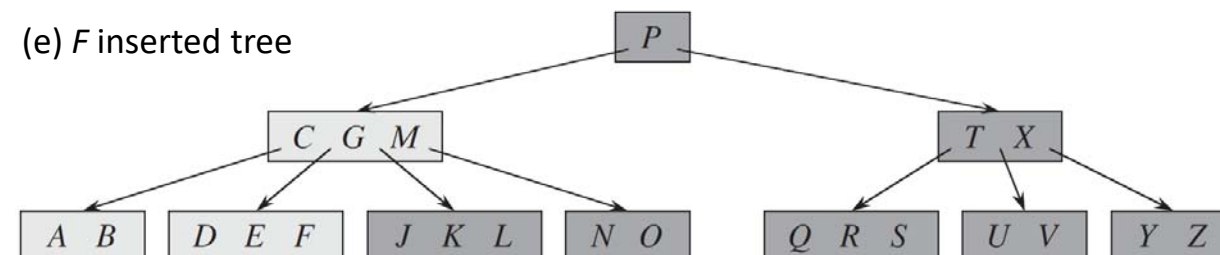
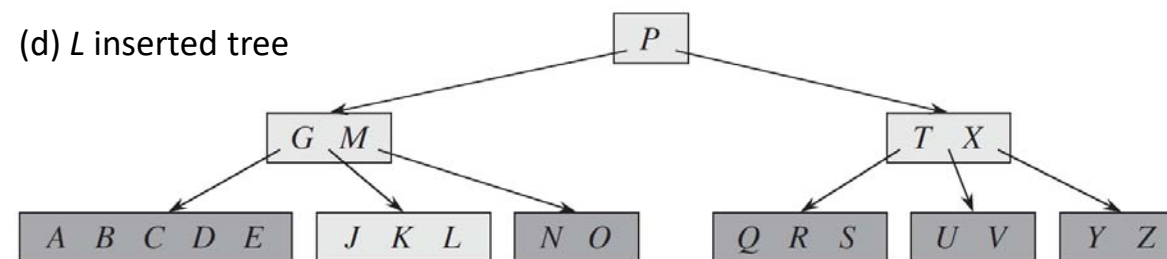
internal: 從右往左比, 找合適的子樹,
在往下之前先檢查兒子
有沒有滿



$t = 3$



leaf: 從右往左比, 找合適的地方放



Outline

- ▶ Definition of B-trees
- ▶ Basic operations on B-trees
- ▶ **Deleting a key from a B-tree 刪除一個鍵值**

從 root 開始往下拜訪, 拜訪的過程中保證拜訪至少要有 t 個 keys

Deleting a key into a B-tree_{1/3} ⇒ 如此可保證只需由上往下一次

- ▶ **B-TREE-Delete** deletes the key k from the subtree rooted at x .
 - ▶ Guarantee that whenever B-TREE-DELETE is called recursively on a node x , the number of keys in x is at least the minimum degree t .
 - ▶ Allows us to delete a key from the tree in one downward pass without having to "back up".

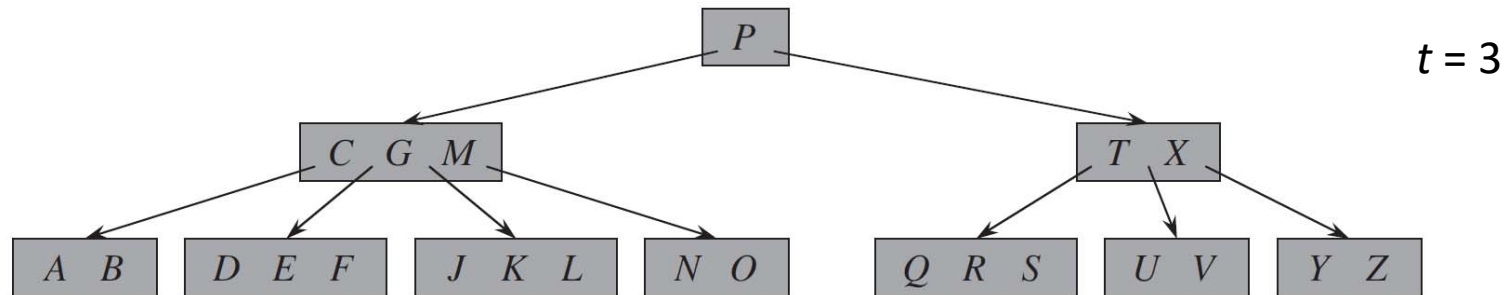
如果 root 沒有 contain 任何鍵值 ⇒ 將 root 刪掉 ⇒ 樹的高度少 1

- ▶ If the root node x becomes an internal node having no keys, then (occur in case 3c, below)
 - ▶ x is deleted,
 - ▶ x 's only child $c_1[x]$ becomes the new root of the tree,
 - ▶ decreasing the height of the tree by one, and
 - ▶ preserving the property that the root of the tree contains at least one key.

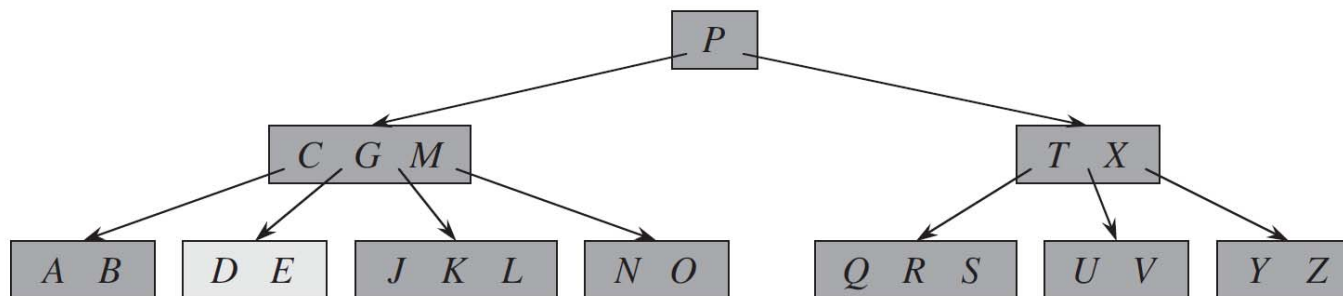
Case 1: k is in node x and x is a leaf

$\left\{ \begin{array}{l} x \text{ 是 leaf, } k \text{ 在 node } x \text{ ①} \\ x \text{ 是 internal node } \left\{ \begin{array}{l} k \text{ 在 node } x \text{ ②} \\ k \text{ 不在 node } x \text{ ③} \end{array} \right. \end{array} \right.$

- ▶ Delete the key k from x .



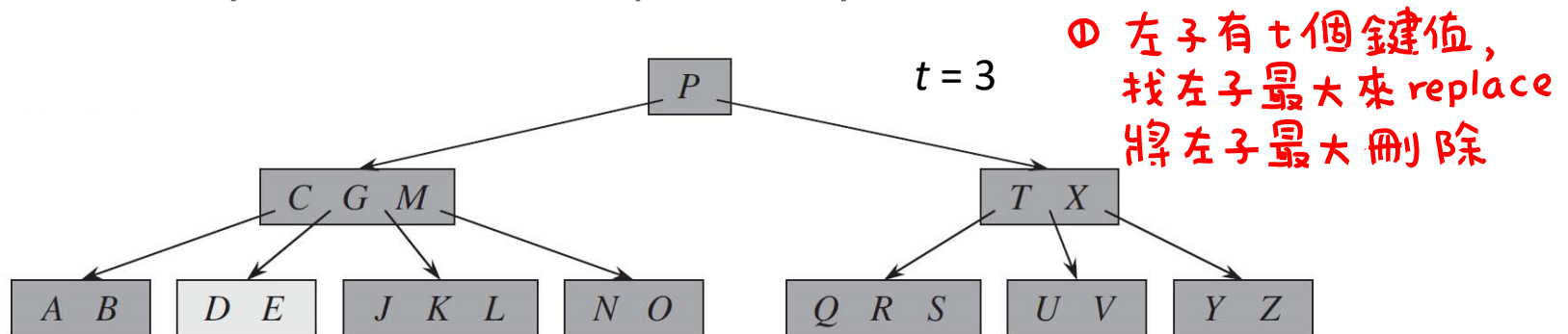
Case 1: F deleted. x 是 leaf, k 在 node $x \Rightarrow$ 直接删



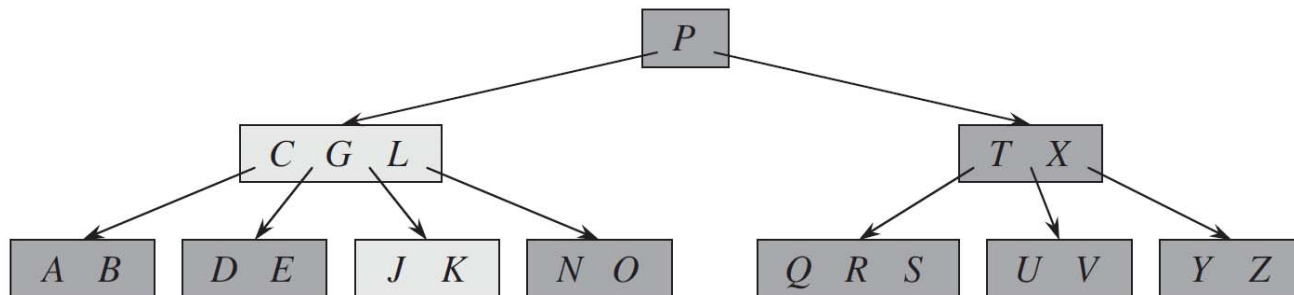
Binary search tree: x 有 2 子 \Rightarrow ① 左子找最大
or ② 右子找最小 來 replace

Case 2: k is in node x and x is an internal node_{1/3}

- ▶ Case 2a: the child y that precedes k has at least t keys.
 - ▶ Find the predecessor k' of k in the subtree rooted at y .
 - ▶ Recursively delete k' , and replace k by k' in x .



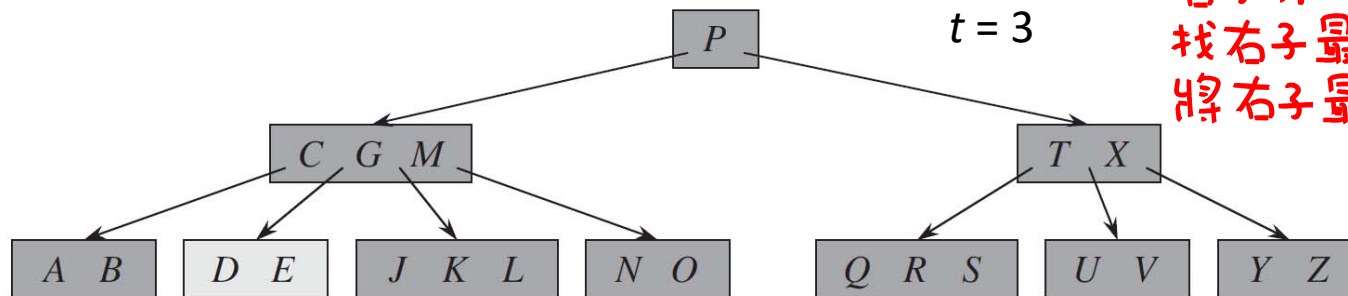
Case 2a: M deleted.



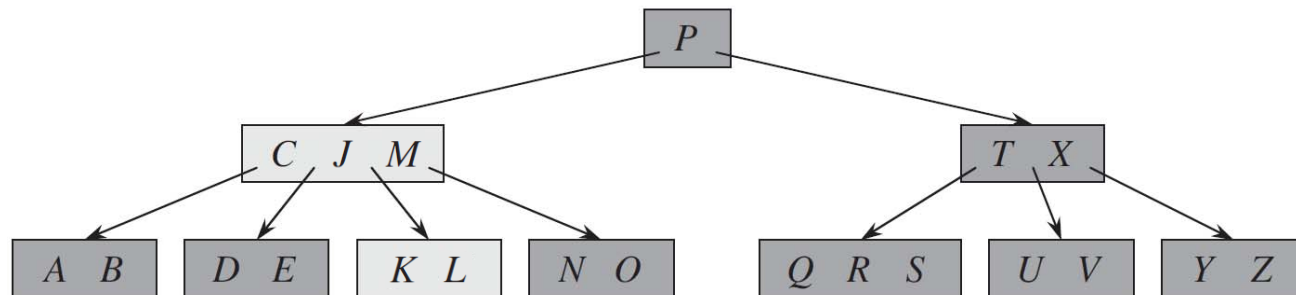
Case 2: k is in node x and x is an internal node_{2/3}

- ▶ Case 2b: the child z that follows k has at least t keys.
 - ▶ Find the successor k' of k in the subtree rooted at z .
 - ▶ Recursively delete k' , and replace k by k' in x .

雖然有 t 個, 但不可以直接刪
② 右子有 t 個鍵值, 找右子最小來 replace 將右子最小刪除



Case 2a: G deleted.



Case 2: k is in node x and x is an internal node_{3/3}

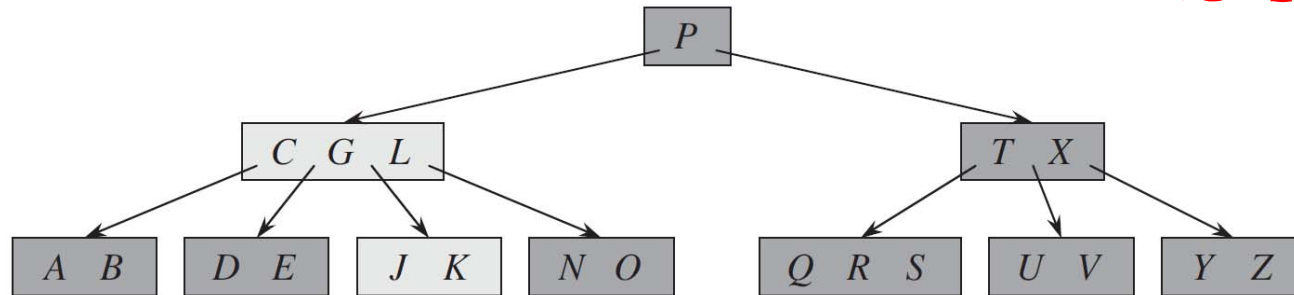
- ▶ Case 2c: both y and z have only $t - 1$ keys. ③ 右子, 左子都只有 $t-1$ 個鍵值
 - ▶ Merge k and all of z into y .
 - ▶ Free z and recursively delete k from y .

⇒ 將右子, 左子, k 合併

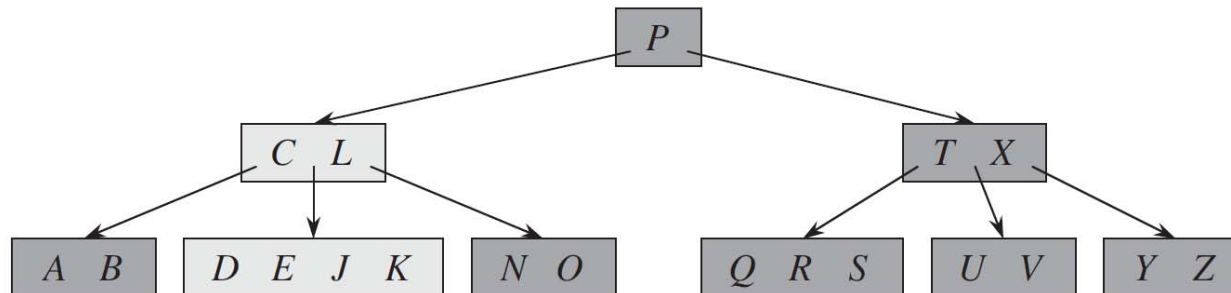
⇒ 用遞迴將鍵值刪除

(已往下-層)

$t = 3$



Case 2c: G deleted.



k 不在 internal node x 中, k 在 $C_i[x]$ 子樹中

Case 3: k is not present in internal node $x_{1/3}$

- ▶ Determine the root $c_i[x]$ of the appropriate subtree that must contain k .

(I) $C_i[x]$ 有 t 個

- ▶ Case 3a: $c_i[x]$ has at least t keys

→ 用遞迴將 k 從 $C_i[x]$

- ▶ Recursively delete k from $c_i[x]$.

子樹中刪除

- ▶ Case 3b: $c_i[x]$ has only $t - 1$ keys but has an immediate sibling with at least t keys.

- ▶ Case 3c: $c_i[x]$ and both of $c_i[x]$'s immediate siblings have $t - 1$ keys.

(II) $C_i[x]$ 只有 $t-1$ 個

→ 兄有 t 個, 跟兄借

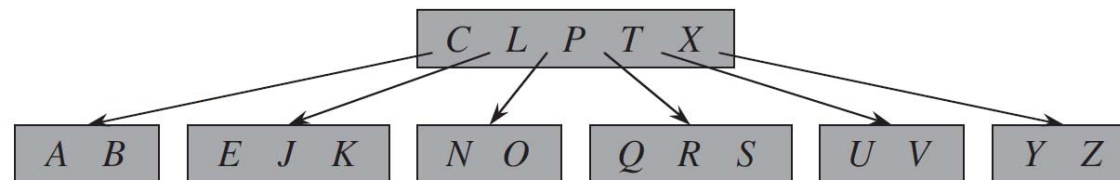
→ 弟有 t 個, 跟弟借

(III) $C_i[x]$, 兄弟都只有 $t-1$ 個

→ 合併

Case 3: k is not present in internal node $x_{2/3}$

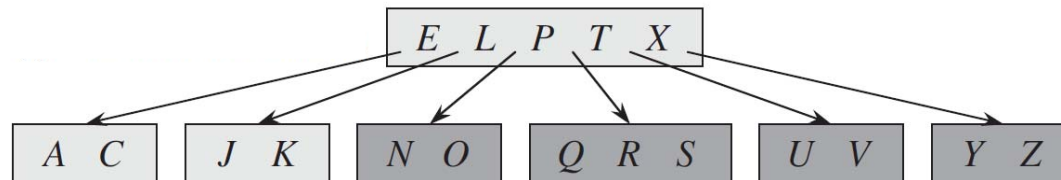
- ▶ Case 3b: $c_i[x]$ has only $t - 1$ keys but has an immediate sibling with at least t keys.
 - ▶ Give $c_i[x]$ an extra key by moving a key from x down into $c_i[x]$.
 - ▶ Moving a key from $c_i[x]$'s immediate left or right sibling up into x .
 - ▶ Moving the appropriate child pointer from the sibling into $c_i[x]$.



$t = 3$

Case 3b : B deleted.

① 兄有 t 個, 跟兄借
② 弟有 t 個, 跟弟借

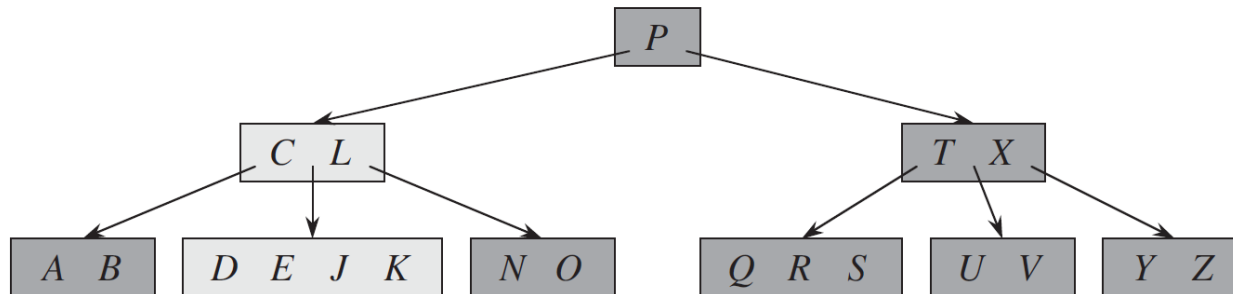


Case 3: k is not present in internal node $x_{3/3}$

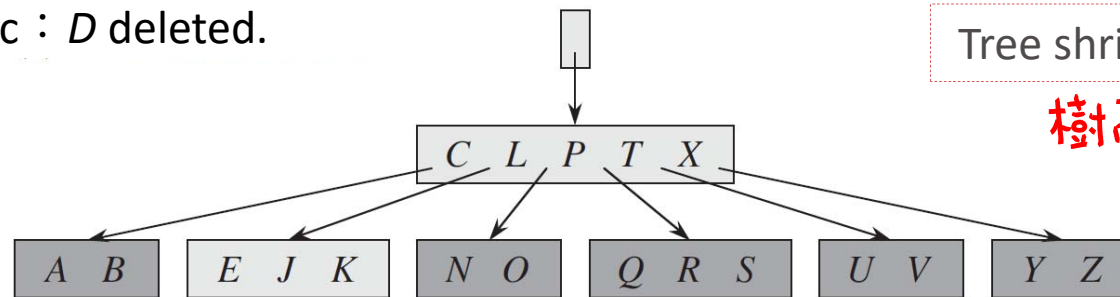
- ▶ Case 3c: $c_i[x]$ and both of $c_i[x]$'s immediate siblings have $t - 1$ keys.
 - ▶ Merge $c_i[x]$ with one sibling.
 - ▶ Moving a key from x down into the new merged node to become the median key for that node.

③ 兄弟都只有 $t-1$ 個, 合併

$t = 3$



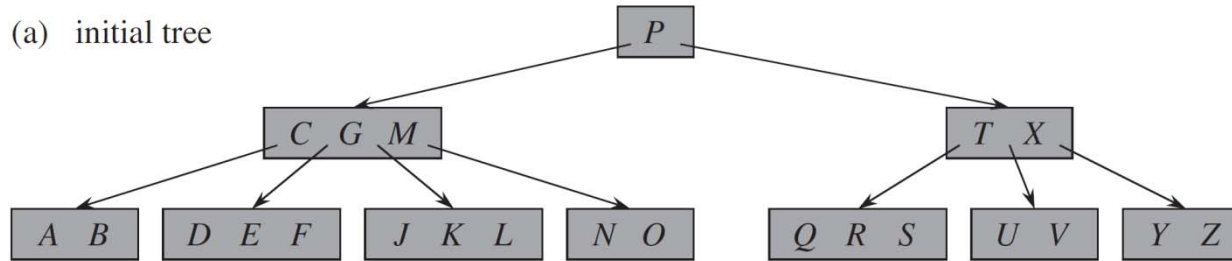
Case 3c : D deleted.



Tree shrinks in height.

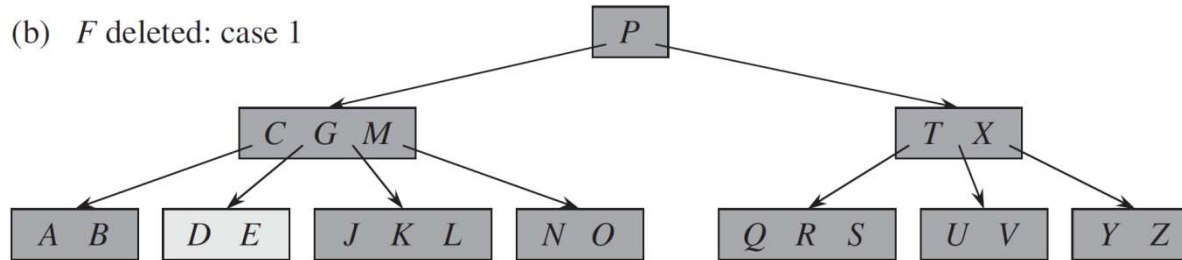
樹高減1

(a) initial tree

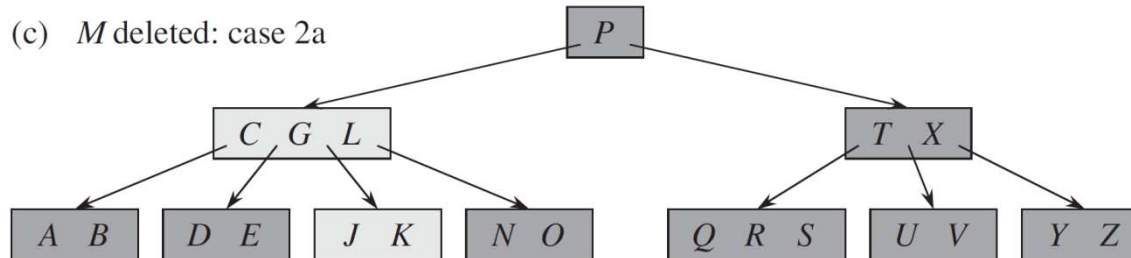


(b) F deleted: case 1

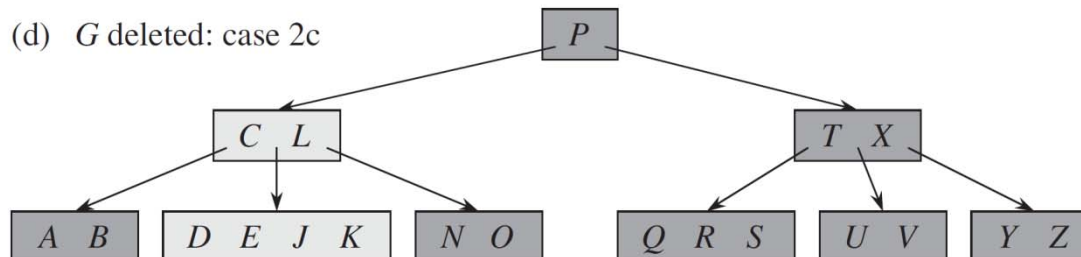
$t = 3$



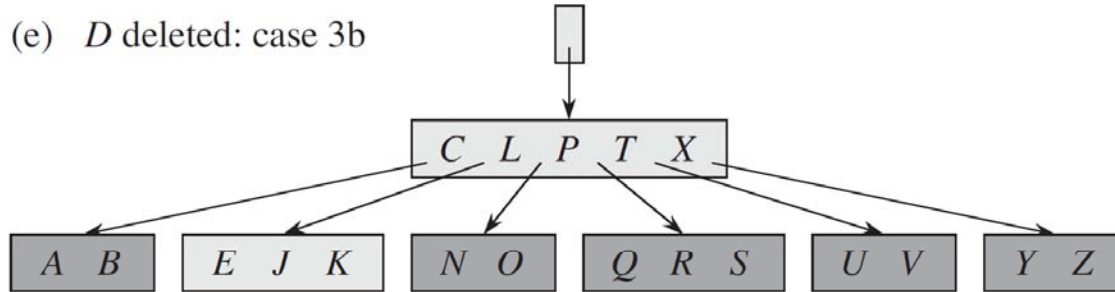
(c) M deleted: case 2a



(d) G deleted: case 2c

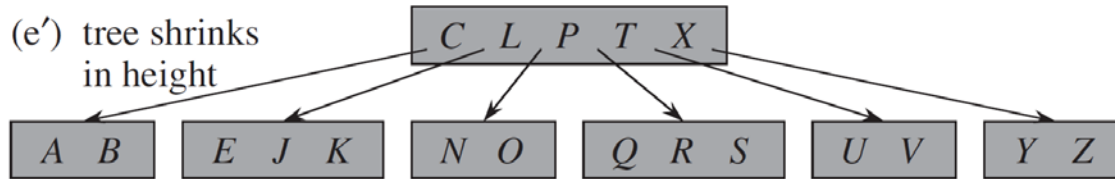


(e) *D* deleted: case 3b

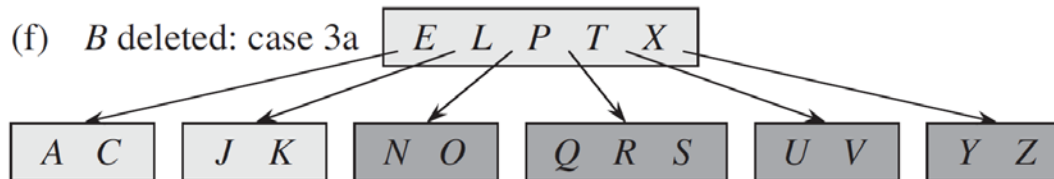


$t = 3$

(e') tree shrinks
in height



(f) *B* deleted: case 3a



Time complexity

- ▶ Only $O(h)$ disk operations for a B-tree of height h .
- ▶ Only $O(1)$ calls to DISK-READ and DISK-WRITE are made between recursive invocations of the procedure.
- ▶ The CPU time required is $O(th) = O(t \log_t n)$.

遞迴程式只有 $O(1)$ 的 DISK-READ 和 DISK-WRITE