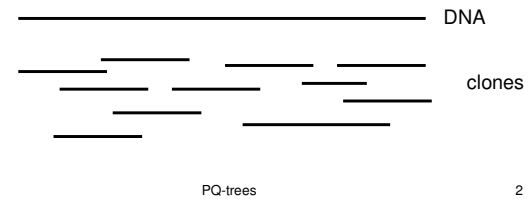


CSEP 521
Applied Algorithms
Autumn 2009

Contiguous Ordering - PQ Trees

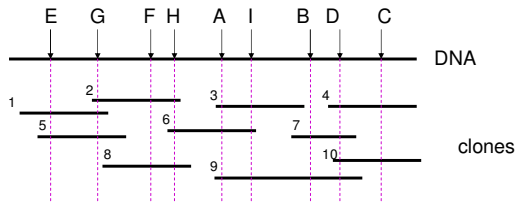
DNA Sequence Reconstruction

- DNA can only be sequenced in relatively small pieces, up to about 1,000 nucleotides.
- By chemistry a much longer DNA sequence can be broken up into overlapping sequences called clones. Clones are 10's of thousands of nucleotides long.



Tagging the Clones

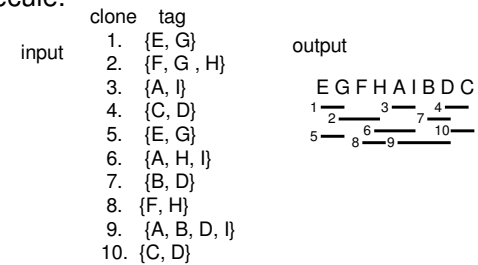
- By chemistry the clones can be tagged by identifying a region of the DNA uniquely.



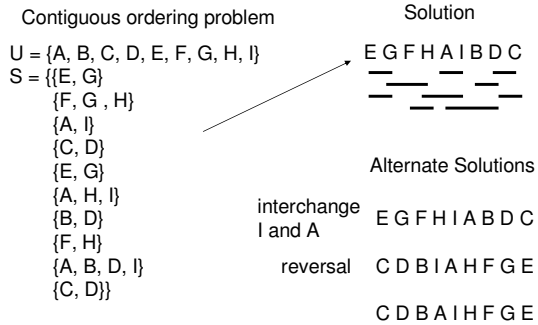
- Each clone is then tagged correspondingly.

Problem to Solve

- Given a set of tagged clones, find a consistent ordering of the tags that determines a possible ordering of the DNA molecule.



Contiguous Ordering Solutions



PQ-trees

5

Linear Time Algorithm

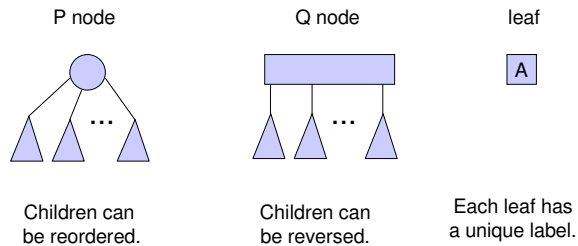
- Booth and Lueker, 1976, designed an algorithm that runs in time $O(n+m+s)$.
 - n is the size of the universe, m is the number of sets, and s is the sum of the sizes of the sets.
- It requires a novel data structure called the PQ tree that represents a set of orderings.
- PQ trees can also be used to test whether an undirected graph is planar.

PQ-trees

6

PQ Trees

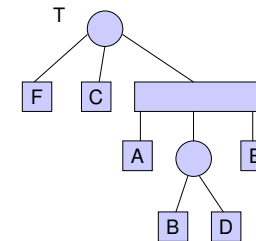
- PQ trees are built from three types of nodes



PQ-trees

7

Example PQ-Tree



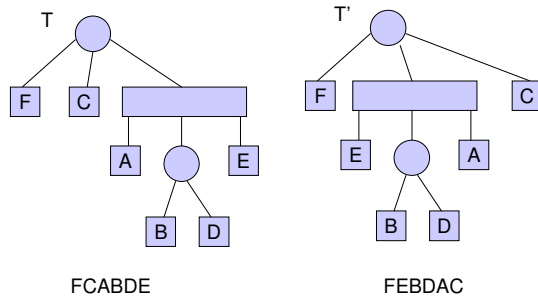
The frontier of T defines the ordering $F(T) = FCABDE$, just read the leaves left to right.

T' is equivalent to T if T' can be transformed into T by reordering the children of P nodes and reversing the children of Q nodes.

PQ-trees

8

Equivalent PQ Trees

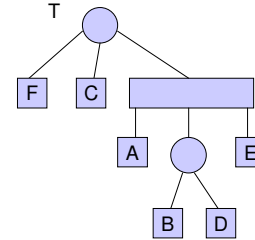


PQ-trees

9

Orderings Defined by a PQ Tree

- Given a PQ tree T the orderings defined by T is
 - $PQ(T) = \{F(T') : T' \text{ is equivalent to } T\}$



There are $6 \times 2 \times 2 = 24$ distinct orderings in $PQ(T)$.

Generally, if a PQ tree T has q Q nodes and p P nodes with number of children c_1, c_2, \dots, c_p , then the number of orderings in $PQ(T)$ is $2^q c_1! c_2! \dots c_p!$.

$$n! = 1 \times 2 \times \dots \times n$$

PQ-trees

10

PQ Tree Solution for the Contiguous Ordering Problem

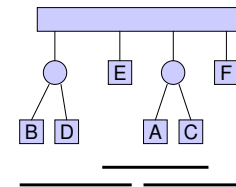
- Input: A universe U and a set $S = \{S_1, S_2, \dots, S_m\}$ of subsets of U.
- Output: A PQ tree T with leaves U with the property that $PQ(T)$ is the set of all orderings of U where each set in S is contiguous in the ordering.

PQ-trees

11

Example Solution

$U = \{A, B, C, D, E, F\}$
 $S = \{\{A, C, E\}, \{A, C, F\}, \{B, D, E\}\}$



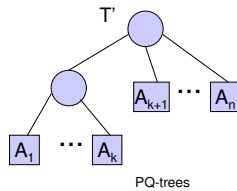
There are 8 orderings that are possible in keeping each of these sets contiguous.

PQ-trees

12

PQ Tree Restriction

- Let $U = \{A_1, A_2, \dots, A_n\}$, $S = \{A_1, A_2, \dots, A_k\}$, and T a PQ tree.
- We will define a function Restrict with the following properties:
 - Restrict(T, S) is a PQ tree.
 - $PQ(\text{Restrict}(T, S)) = PQ(T) \text{ intersect } PQ(S)$ where



13

High Level PQ tree Algorithm

- Input is $U = \{A_1, A_2, \dots, A_n\}$, and subsets S_1, S_2, \dots, S_m of U .
- Initialization:
 - $T = P$ node with children A_1, A_2, \dots, A_n
- Calculate m restrictions:
 - for $j = 1$ to m do
 - $T := \text{Restrict}(T, S_j)$
- At the end of iteration k :
 - $PQ(T)$ = the set of ordering of U where each set S_1, S_2, \dots, S_k are contiguous.

PQ-trees

14

Marking Nodes

- Given a set S and PQ tree T we can mark nodes either full or partial.
 - A leaf is full if it is a member of S .
 - A node is full if all its children are full.
 - A node is partial if either it has both full and non-full children or it has a partial child.
 - A node is doubly partial if it has two partial children.

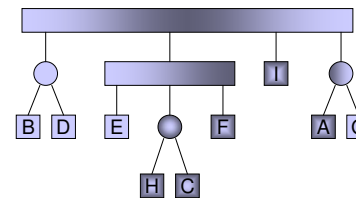
PQ-trees

15

Marks of Nodes

Mark the leaves in S full.
Bottom up mark the nodes full or partial.
The members of S will become contiguous.

$S = \{A, C, F, H, I\}$



Full

left partial

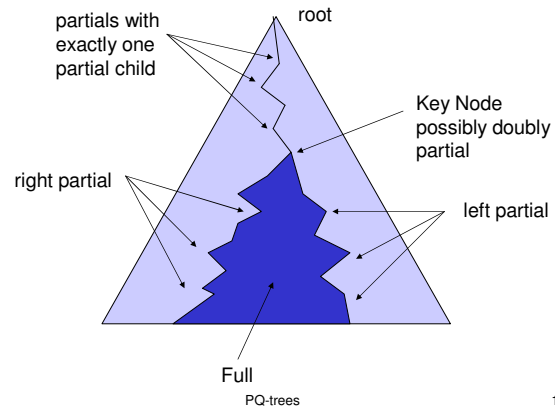
right partial

doubly partial

PQ-trees

16

Structure of the Marked PQ Tree



17

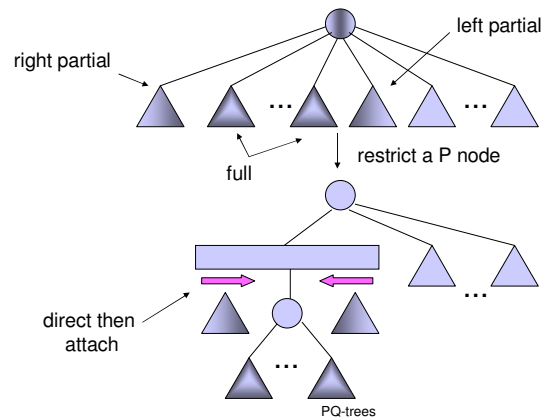
Restrict(T,S)

- Mark the full and partial nodes from the bottom up.
 - In the process the marked leaves become contiguous.
- Locate the key node.
 - Deepest node with the property that all the full leaves are descendants of the node.
- Restrict the key node.
 - In the process of restricting the key node we will have to recursively direct partial nodes.
 - Directing a node returns a sequence of nodes.

PQ-trees

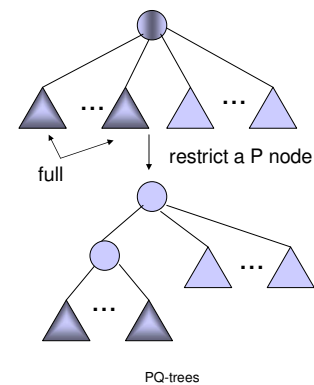
18

Restricting a P Node with Partial Children



19

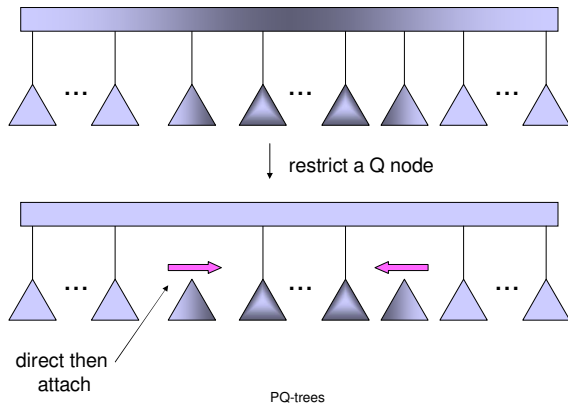
Restricting a P node with no Partial Children



PQ-trees

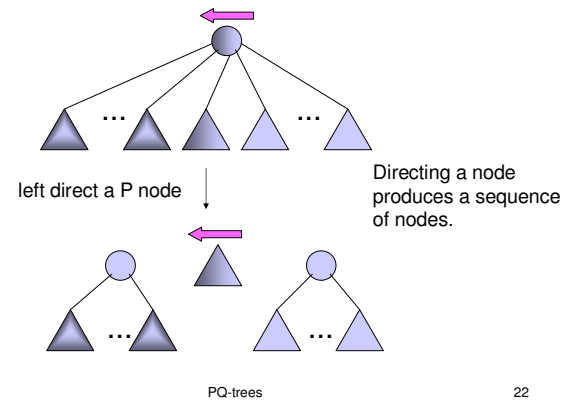
20

Restricting a Q node



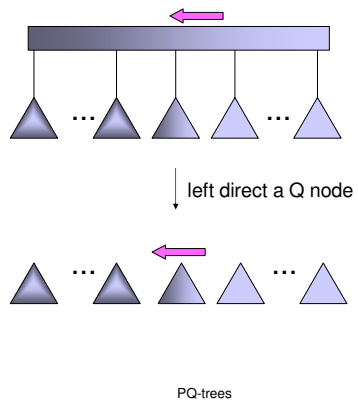
21

Directing a P Node



22

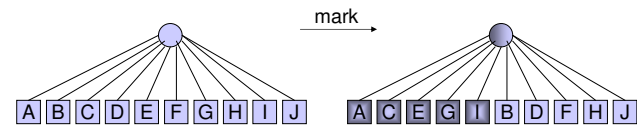
Directing a Q Node



23

Example (1)

$U = \{A, B, C, D, E, F, G, H, I, J\}$
 $S_1 = \{A, C, E, G, I\}$

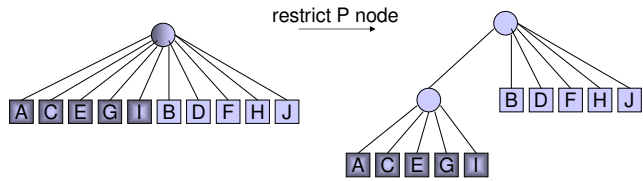


PQ-trees

24

Example (2)

$U = \{A, B, C, D, E, F, G, H, I, J\}$
 $S_1 = \{A, C, E, G, I\}$



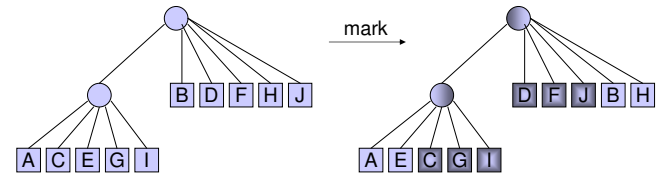
special case because
no partial child.

PQ-trees

25

Example (3)

$U = \{A, B, C, D, E, F, G, H, I, J\}$
 $S_2 = \{C, D, F, G, I, J\}$

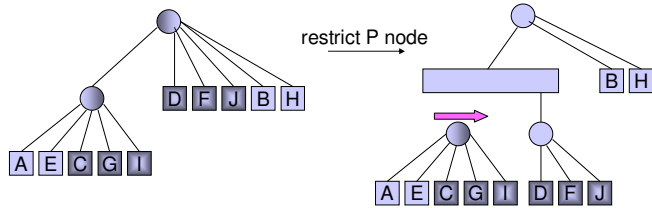


PQ-trees

26

Example (4)

$U = \{A, B, C, D, E, F, G, H, I, J\}$
 $S_2 = \{C, D, F, G, I, J\}$

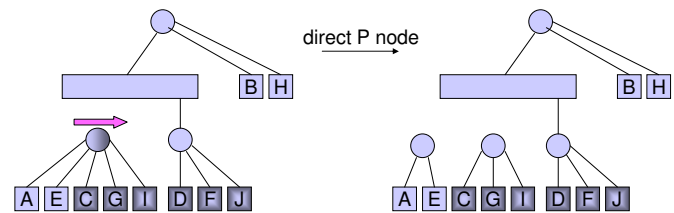


PQ-trees

27

Example (5)

$U = \{A, B, C, D, E, F, G, H, I, J\}$
 $S_2 = \{C, D, F, G, I, J\}$

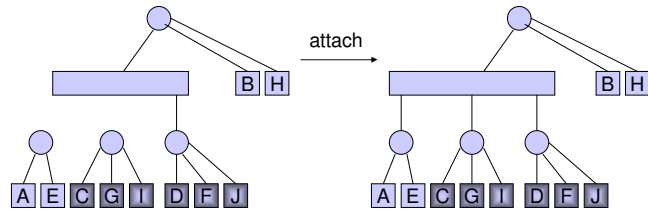


PQ-trees

28

Example (6)

$U = \{A, B, C, D, E, F, G, H, I, J\}$
 $S_2 = \{C, D, F, G, I, J\}$

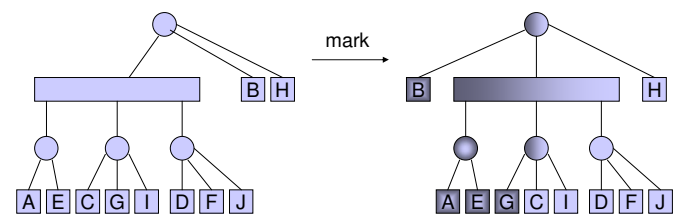


PQ-trees

29

Example (7)

$U = \{A, B, C, D, E, F, G, H, I, J\}$
 $S_3 = \{A, B, E, G\}$

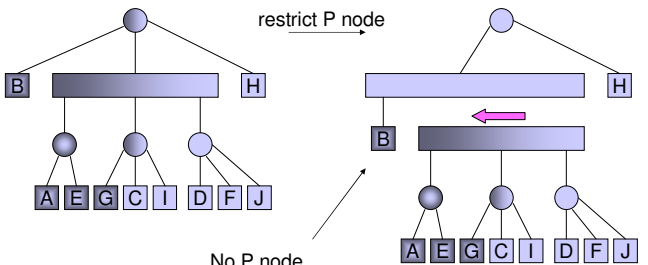


PQ-trees

30

Example (8)

$U = \{A, B, C, D, E, F, G, H, I, J\}$
 $S_3 = \{A, B, E, G\}$

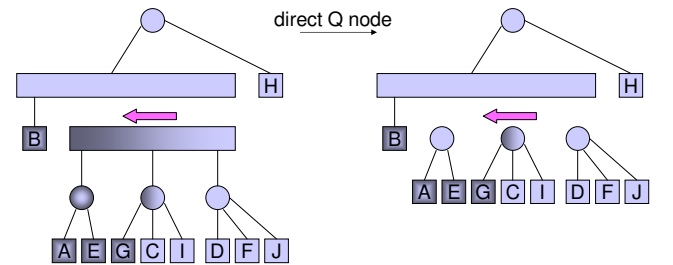


PQ-trees

31

Example (9)

$U = \{A, B, C, D, E, F, G, H, I, J\}$
 $S_3 = \{A, B, E, G\}$

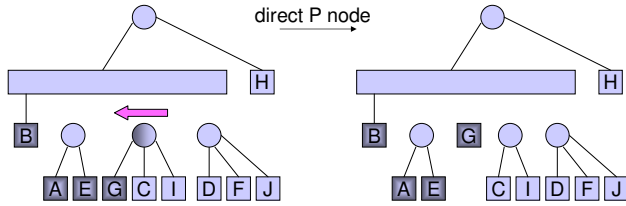


PQ-trees

32

Example (10)

U = {A,B,C,D,E,F,G,H,I,J}
 S₃ = {A,B,E,G}

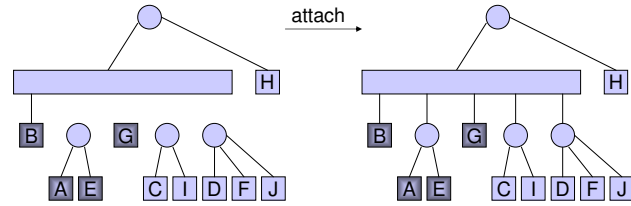


PQ-trees

33

Example (11)

U = {A,B,C,D,E,F,G,H,I,J}
 S₃ = {A,B,E,G}

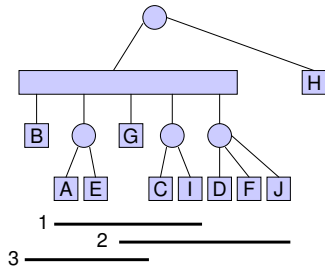


PQ-trees

34

Example (12)

U = {A,B,C,D,E,F,G,H,I,J}
 S₁ = {A,C,E,G,I}
 S₂ = {C,D,F,G,I,J}
 S₃ = {A,B,E,G}

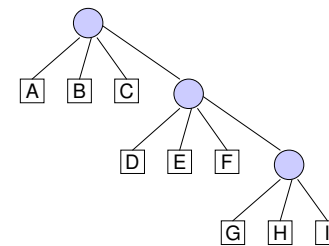


PQ-trees

35

Exercise

- Restrict with to make {A,B,D,E,G} contiguous



PQ-trees

36

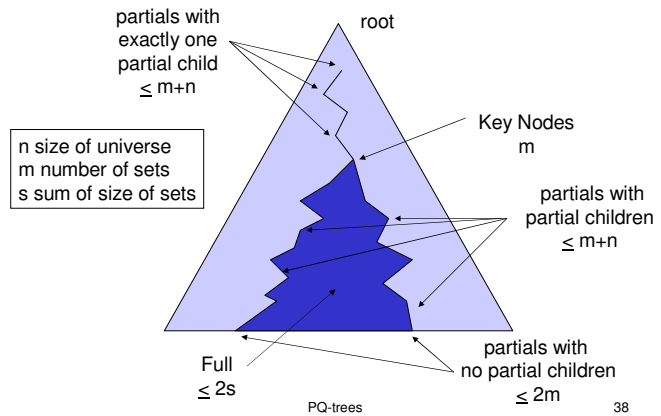
Linear Number of Nodes Processed

- Let n be the size of the universe, m the number of sets, and s the sum of the sizes of the sets.
 - Number of full nodes processed $\leq 2s$.
 - Number of key nodes processed $= m$.
 - Number of partial nodes with partial children processed below the key node $\leq m + n$.
 - Number of partial nodes with no partial children $\leq 2m$.
 - Number of partial nodes processed above the key node $\leq m + n$.

PQ-trees

37

Number of Processed Nodes Amortized



38

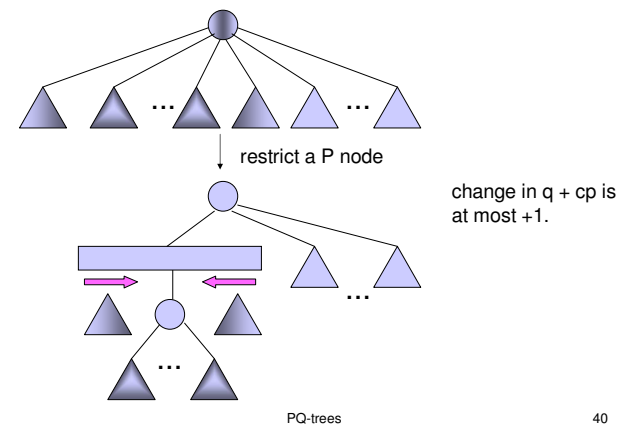
Partials with Partial Children Below the Key Node

- Amortized complexity argument.
- Consider the quantities:
 - q = number of Q nodes,
 - cp = number of children of P nodes.
 - We examine the quantity $x = q + cp$
 - x is initially n and never negative.
 - Each restrict of a key node increases x by at most 1.
 - Each direct of a partial node with a partial child decreases x by at least 1.
 - Since there are m restricts of a key node then there are most $n + m$ directs of partials with partial children.

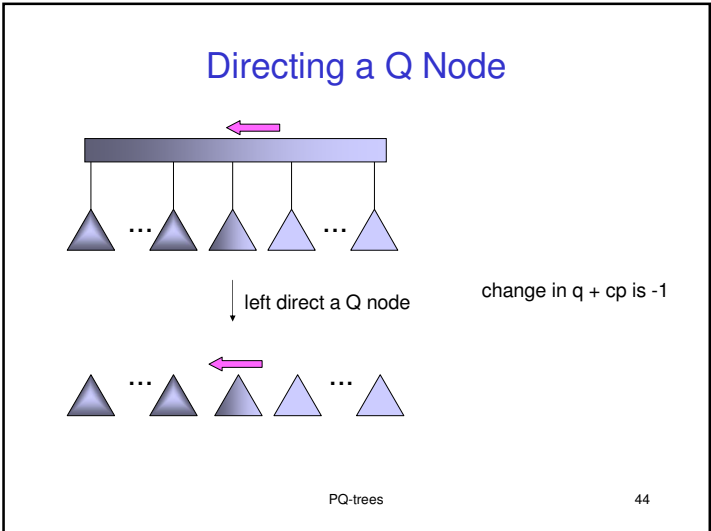
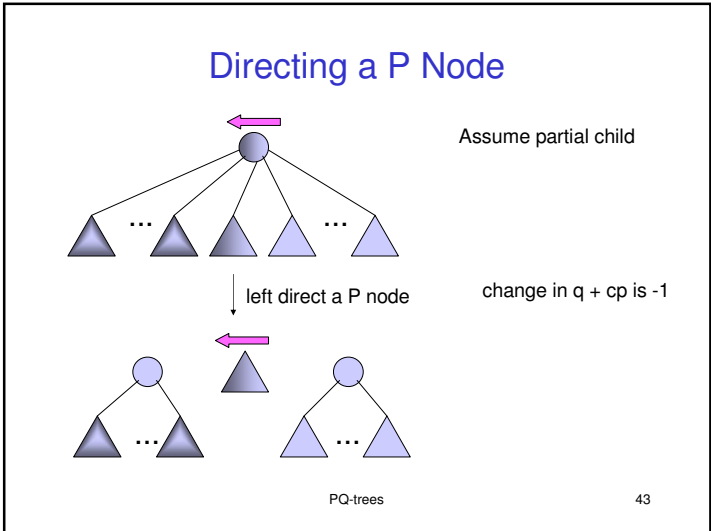
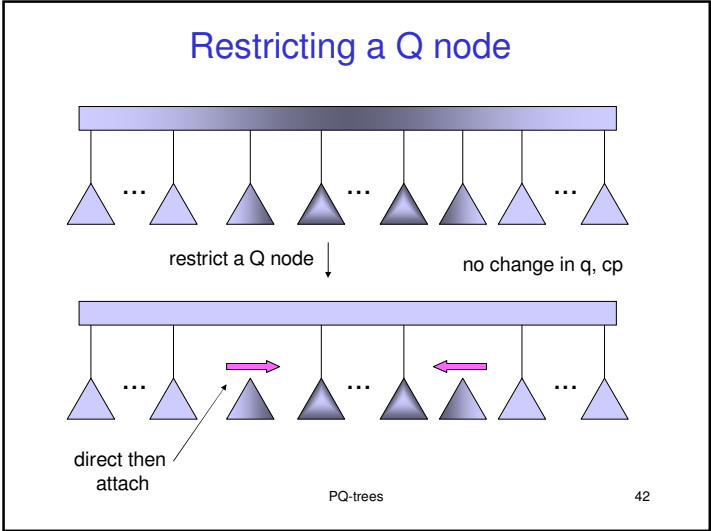
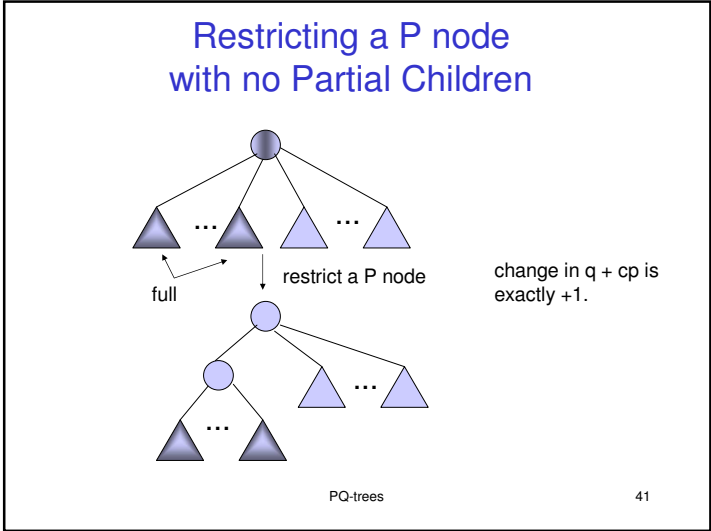
PQ-trees

39

Restricting a P Node with Partial Children



40



PQ Tree Notes

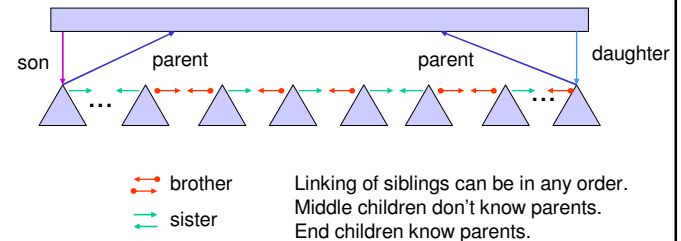
- In algorithmic design only a linear number of nodes are ever processed.
- Designing the data structures to make the linear time processing a reality is very tricky.
- PQ trees solve the idealized DNA ordering problem.
- In reality, because of errors, the DNA ordering problem is NP-hard and other techniques are used.

PQ-trees

45

Example of Data Structure Trick

- Linking the children of a Q node



PQ-trees

46