

XML Query Evaluation

* XML Fragmentation:

=> XML Fragmentation refers to the process of breaking down large XML Document, into smaller.

Fragmentation is a Process of breaking or being broken into Fragment.

Fragmentation can improve the scalability of the XML Document because we can easily share and store the whole XML Document into the small Fragment.

There are Two Types of Fragmentation

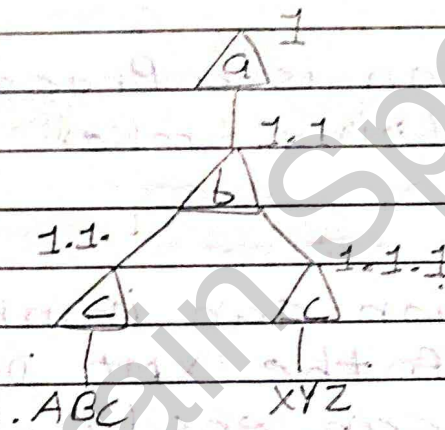
1) Horizontal: Fragmentation based on selection operators and predicates.

2) Vertical: Fragmentation based on a partitioning of the set of elements type in a Schema.

Ex. `<a>`
``
`<c> ABC </c>`
`<d> XYZ </d>`
``
``

XML Data

=> XML Fragment tree with XFL



=> XML Fragment with XFL

`<Fragment FID="1" tsid="11">`
`<a/>``</Fragment>`

`<Fragment FID="1.1" tsid="12">`
```<c> ABC </c>```
`</Fragment>`

`<Fragment FID="1.1.1" tsid="13">`
`<c> XYZ </c>`
`</Fragment>`

* XML Identifier : Region-Based Identifier

⇒ XML Region-Based Identifier are a way to identify and manage the structure of XML document.

In XML Region-Based Identifier use offset or positional information to identify XML document.

This method use pairs of position to denote the start and end of the each element.

→ Basic of Region-Based Identifier or Steps of Region-Based Identifier

1 (Begin Tag, End Tag) Counting :

Instead of counting character offsets, count only opening and closing tag.

Count offset start and ~~Eng~~ End.

2 (Pre, Post) Identifiers :

Count nodes during a preorder traversal and assign these as <pre> value.

Count nodes during a post-order traversal and assign this as 'post' value.

3 (pre, post, Depth) Identifiers:

Adds a 'depth' value, which representing the distance from the root.

=> IF 'n1' is ancestor of 'n2'
 iff $n1.pre \leq n2.pre$ and
 $n2.post \leq n1.post$

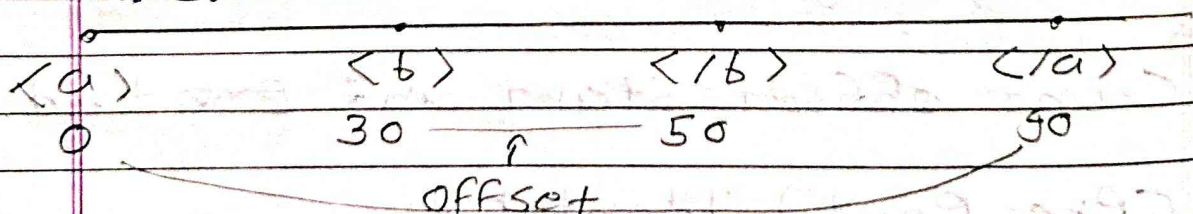
=> Node 'n1' is parent of n2
 then $n1.depth = n2.depth - 1$

Ex. <a>

Some Text

 Text

More



=> Here, For element

<a> offset = (0, 90)

 offset = (30, 50)

\Rightarrow Here, $\langle a \rangle$ is a parent of b .
So, Starting Tag $\langle a \rangle$ and Ending
Tag $\langle /a \rangle$

So, Tag Count = (1, 4)

For, $\langle b \rangle$, Tag Count = (2, 3)

\Rightarrow Here, $\langle a \rangle$ is Root element and
parent of $\langle b \rangle$ (Exchange)

So, ~~at~~ $a.pre = b.pre$
1 = 2

So, $a.post = 2$ and
 $b.post = 1$

\Rightarrow Here, $\langle a \rangle$ is Root element
So, distance from Root = Depth = 0

and For $\langle b \rangle = 1$

So, For $\langle a \rangle = (pre, post, depth)$
= (1, 2, 0)

$\langle b \rangle = (pre, post, depth)$
= (2, 1, 1)

* XML Identifier: Dewey Classification

=> The Dewey classification schema, originally designed for library organization and can be used for XML Document.

Dewey is provide hierarchical structure or Identifier system to the XML Document.

This method creates identifiers by appending suffixes to the parent's ID.

=> Principle of Dewey IDs

1 Hierarchical Structure:

Each node's ID is derived by appending a unique suffix to its parent's ID.

IF Node n_1 is parent of n_2 than n_1 is a prefix of n_2

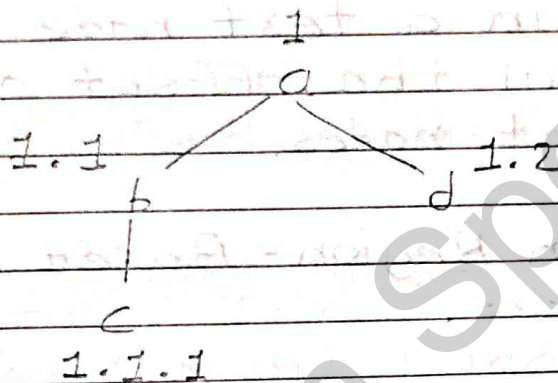
Ex. n_1 - prefix = 1.1 than
 $n_2 = 1.1.1$ (extends)

IF n_2 has sibling called n_3 than $n_3 = 1.1.2$ (same Parent Prefix)

Ex. <a>

 <c> Text </c>

 <d> More </d>



* Structural Identifiers and Updates in XML Document

=> When XML Document are updates or change than structural identifier efficiently maintaining the XML Document's Integrity and performance.

Structural Identifiers Manage the changes of XML Document Structure

=> Impact :

1 Change in offset-Based Identifiers:

Even a single character add or remove in a text node can change all the offset of all subsequent nodes.

2 Change in Region-Based Identifiers.

This Identifier are not impact by the changes to text node.

Inserting node require adjusting the identifier of all subsequent node in pre order.

Deleting node require to introduce gap, but not impact on structural join.

3 Dewey IDs:

Character change / remove can not impacted to change text nodes.

Inserting a new node require the change all subsequent nodes Dewey ID's

Deletion can gives a gap, but not required Dewey ID changes

* XML Evaluation Techniques : Structural Join.

=> Structural Join in XML databases are fundamental operation that combine tuples from two input based on a structural relationship

=> Let P_1 and P_2 be partial evaluation plans in XML Database, where

- $P_1.x$ is an attribute of output of P_1
- $P_2.y$ is an attribute of output P_2

Both, $P_1.x$ and $P_2.y$ has structural IDs.

-> Two Key binary relationship are used,

- is parent of (\prec): One node is parent of another.

- isAncestorof(\ll): One node is an Ancestor of another.

Structural Join,

- IF $P_1.X$ is ancestor of $P_2.Y$

$$P_1 \bowtie X \ll Y P_2 = \{ (t_1, t_2) \mid t_1 \in P_1, t_2 \in P_2, t_1.X \ll t_2.Y \}$$

- IF $P_1.X$ is parent of $P_2.Y$

$$P_1 \bowtie X < Y P_2 = \{ (t_1, t_2) \mid t_1 \in P_1, t_2 \in P_2, t_1.X < t_2.Y \}$$

=> Efficiency Consideration:

Goal is to reduce CPU, memory and I/O cost while we perform structural join on large XML Document.

1 Nested Loop Join: Iterates over the output of P_1 for each tuple, iterates over the output of P_2

$$\text{Cost} = O(|P_1| \times |P_2|)$$

2 Hash Join: Similar to Hash Join in Relational Database

$$\text{Cost} = O(|P_1| + |P_2|)$$

3 stack-Based Join: Used to evaluate scheme that have (start, end) ID

$$\text{Cost} = \Omega(|P_1| + |P_2|)$$

Ex. <book>

<book id="1">

<+> XYZ </+>

<a> A

</book>

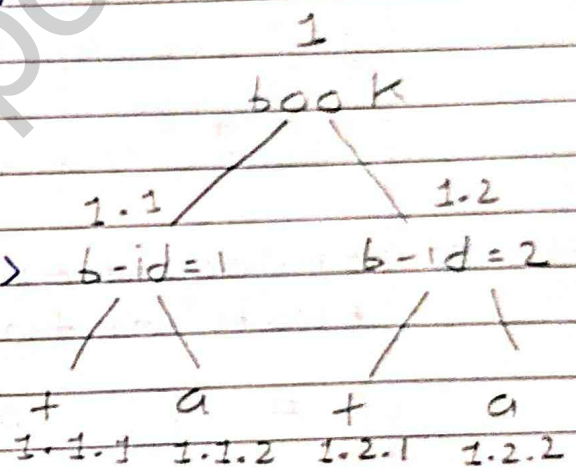
<book id="2">

<+> ABC </+>

<a> X

</book>

</book>



We want to use Stack to Join <book> and <+> element.

Stack 1 (P₁) = Hold <book> element

Stack 2 (P₂) = Hold <+> element

Stack 3 (P₃) = Output

→ Steps :

1 Initial state: Three Stack are empty.

2 Encounter $\langle \text{book id} = "1" \rangle$

Push onto stack 1: P_1 ($b-id=1$)

P_1	P_2	O/P
1.1		

3 Encounter $\langle + \rangle$ & $b-id=1$

Push onto stack 2: P_2 ($1.1.1$)

P_1	P_2	O/P
1.1	1.1.1	

Condition Check:

$P_1.top() < P_2.top()$

1.1 is parent of 1.1.1

So, This is valid Pair

So, Push into O/P stack

P_1	P_2	O/P
1.1	1.1.1	$(1.1), (1.1.1)$

4 Same For <book id = "2">

Push onto stack P₁ (b-id=2)

Push onto stack P₂ @ b-id=2 → <+>

Condition check:

$P_1.top() < P_2.top()$

1.2 is a parent of 1.2.1

So, Push into O/P stack.

P ₁	P ₂	O/P
1.2	1.2.1	
1.1	1.1.1	(1.1), (1.1.1)

↓ Valid condition

P ₁	P ₂	O/P
1.2	1.2.1	(1.2), (1.2.1)
1.1	1.1.1	(1.1), (1.1.1)

* Optimizing Structural Join:

⇒ There are Two Algorithms are used for Structural Join.

1) STD Algorithm: Combines two input based on matching ID's of node.

2) STA Algorithm: Similar to STD but may sort based on different JD's depending on join order.

For Optimizing the Structural Join We have Four Plan,

1 Plan (a): Logical Plan:

Illustrates a tree pattern query and the initial logical plan evaluating it.

2 Plan (b): Implement Logical plan using STD Algorithm for First Join.

A Sort Operator is required to ensure the correct order of third Join.

3 Plan (c): Uses STA for First Plan and Required Two Sort operation,

One After First Join and another after second Join

4 PlanCDJ: Use STA For both Join, does not require any sort operators.

5 PlanCDJ: Shows a Fully pipelined plan using TwigStack algorithm.

* Holistic Twig Joins:

=> Holistic Twing Joins represents an advanced method to process the entire query.

It is reducing both running time and space requirement.

-> Logical Defination:

For a Query Tree Pattern 'q' with 'n' nodes, where each node 'i' is labeled 'a-i' and 'a-1' is the Root.

The Logical Holistic Twing operator denoted as $\cdot lq C1, p1, 1, p2, \dots, 1, pn$

where, $1, p1, \dots, 1, pn$ are Logical sub-plans the output structure ID.

The result is

$$\sigma(a_{p_2} \ll a_2 \wedge a_{p_3} \ll a_3 \wedge \dots \wedge a_{p_n} \ll a_n) (1_{p_1} \times 1_{p_2} \times \dots \times 1_{p_n})$$

-> Algorithm PathStack

1 Initialization:

Maintain one stack ('S_i') for each query node labeled 'a_i'

Push node ID's onto their respective stack.

2 Execution:

Identify the smallest document order ID.

Push the corresponding ID onto its stack and update pointer between stacks to reflect parent-child relationship

3 Pop ID's from stack that cannot contribute to result tuple

3 Result Construction:

When an entry is pushed onto the stack of a leaf node, build result from entries in all stacks.

* Automata on Ranked Tree:

=> In a Rank tree, the number of a children node can have is Fixed based on the node's label

Each type of node is associated with a specific arity or a Fixed number of child node.

In a binary tree, every internal node has exactly two children.

Automata designed for ranked trees rely on his Fixed structure to define transition and states.

Each node has a Fixed, predefined number of children and gives a uniform and predictable structure.

=> Bottom-Up Tree Automata:

Bottom-Up Tree binary tree start the process from the leaves towards the root.

- Leaf Alphabet: A Finite set of symbols for the leaves
- Internal Alphabet (Σ): A Finite set of symbols for the internal.
- Start Set (Q): A Finite set of states
- Accepting States (F): A subset of states considered as accepting.
- Transition Function (δ): Describes how states are assigned to nodes based on their labels and the states of their children.

=> Top-Down Tree Automata:

Top-Down Tree Automata process the tree from the root towards the leaves.

* Unranked Tree:

In a Unranked tree, there is no Fixed limit on the number of children a node can have.

A node can have any number of children.

In an XML structure, elements can have varying numbers of child elements.

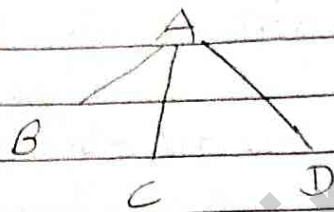
=> Automata on Unranked Tree:

Defining automata for unranked tree is more complex, because the automata needs to handle nodes with varying numbers of children.

Nodes can have a variable number of children, allowing for a more flexible and adaptable tree structure.

Ex. Ranked Tree:

Ex. UnRanked Tree:



=> Ranked Tree:

