

Processing Lax XML Element Trees

Phil Fearon – DeltaXML Ltd (Senior Software Engineer)

Gursheen Kaur – DeltaXML Ltd (Software Developer)

Lax vs Strict – From the Processing Perspective

Lax Content Model



- Many conditional branches
- Complex processing logic
- Reduced opportunity for code reuse
- Comparison of is difficult

Strict Content Model



- Consistent input structure
- Improved processing performance
- Simpler, more robust, code
- Improved test coverage

Lax to Strict – HTML Tables

3

1. *tr* element is a child of *tbody*
2. *tr* element is a child of *tr*
3. two *td* elements follow the *tr*

Lax

```
<table>
  <tbody>
    <tr>
      <td>Anna</td>
      <td>anna@gmail.com</td>
    <tr>
      <td>Bill</td>
      <td>bill@gmail.com</td>
    <td>Charlie</td>
    <td>charlie@gmail.com</td>
  </tbody>
</table>
```

Valid

Hoist

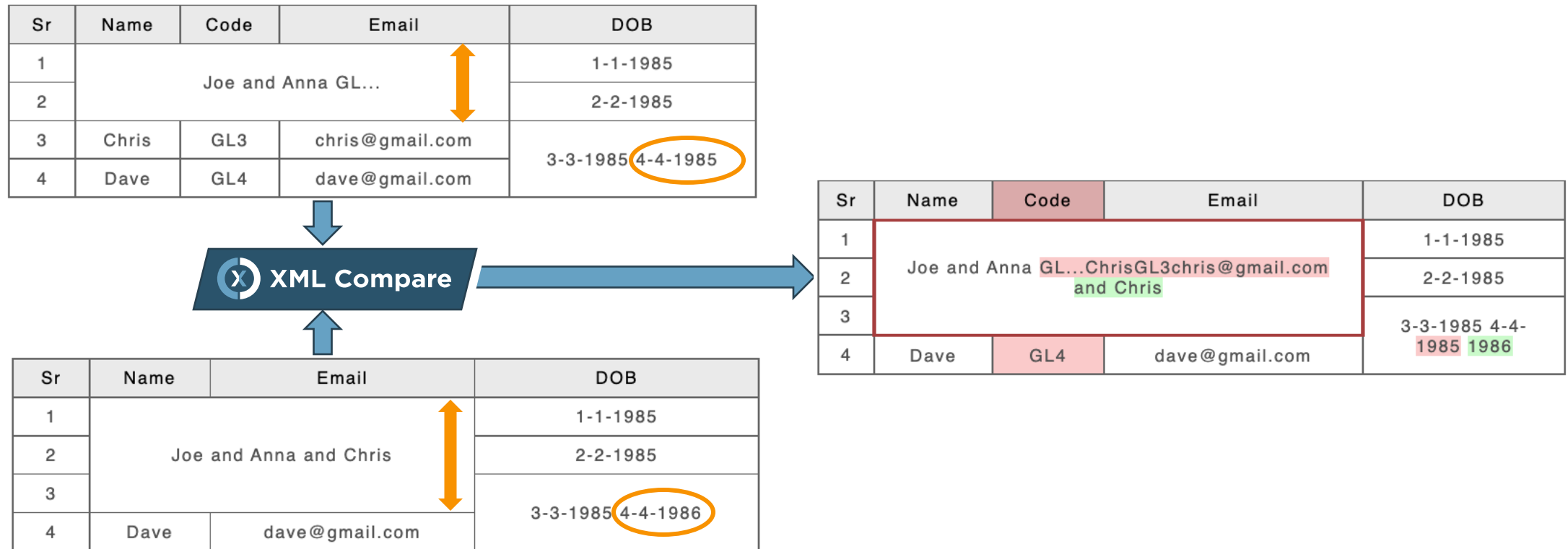
Hoist & Wrap

Strict

```
<table>
  <tbody>
    <tr>
      <td>Anna</td>
      <td>anna@gmail.com</td>
    </tr>
    <tr>
      <td>Bill</td>
      <td>bill@gmail.com</td>
    </tr>
    <tr>
      <td>Charlie</td>
      <td>charlie@gmail.com</td>
    </tr>
  </tbody>
</table>
```

Goal: Compare Lax HTML Tables

With different element structures AND different table-cell spans



HTML Table 'Normalizer' Considerations

- Equivalence with HTML Standard's Table Processing Model
 - Positioning within XSLT pipeline
 - Additional Requirements
 - Relative Performance
 - Non-Requirements
-
- Off-the-shelf HTML parser?
 - How to Test and Verify?

XSLT: Equivalence with HTML Table Processing Model

HTML

Living Standard — Last Updated 30 June 2023

[13 The HTML syntax](#) — [Table of Contents](#)

13.2 Parsing HTML documents

13.2.6 Tree construction

13.2.6.4 The rules for parsing tokens in HTML content

13.2.6.4.9 The "in table" insertion mode

13.2.6.4.10 The "in table text" insertion mode

13.2.6.4.11 The "in caption" insertion mode

13.2.6.4.12 The "in column group" insertion mode

13.2.6.4.13 The "in table body" insertion mode

13.2.6.4.14 The "in row" insertion mode

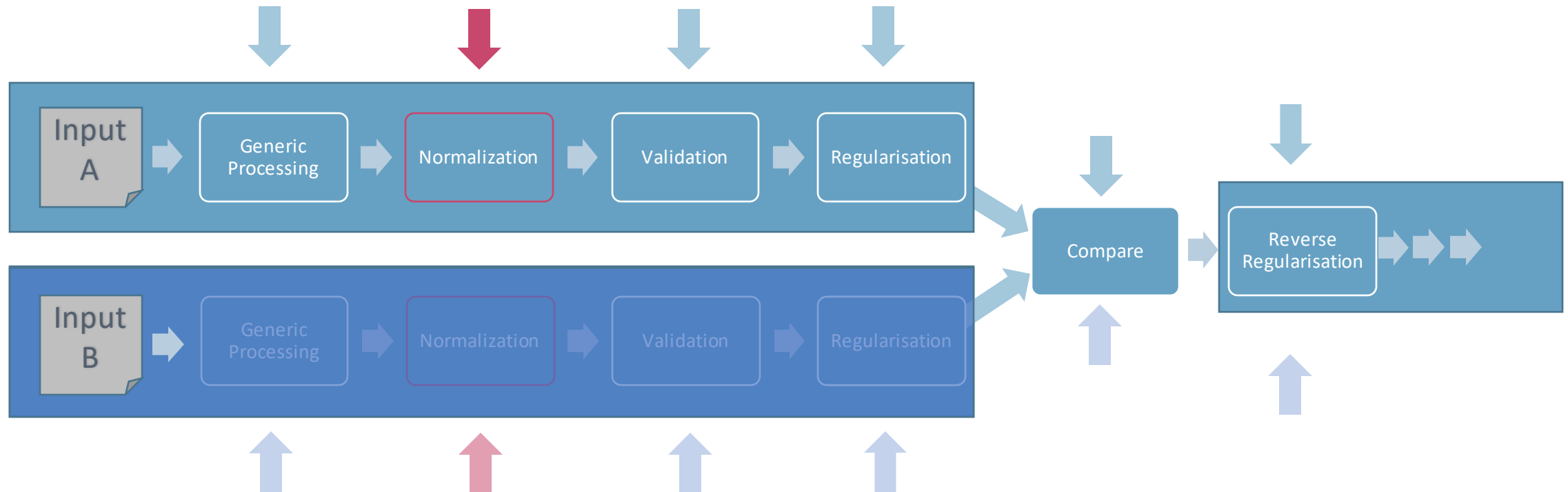
13.2.6.4.15 The "in cell" insertion mode



XSL Transformations (XSLT) Version 3.0

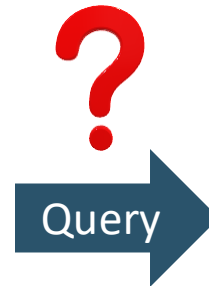
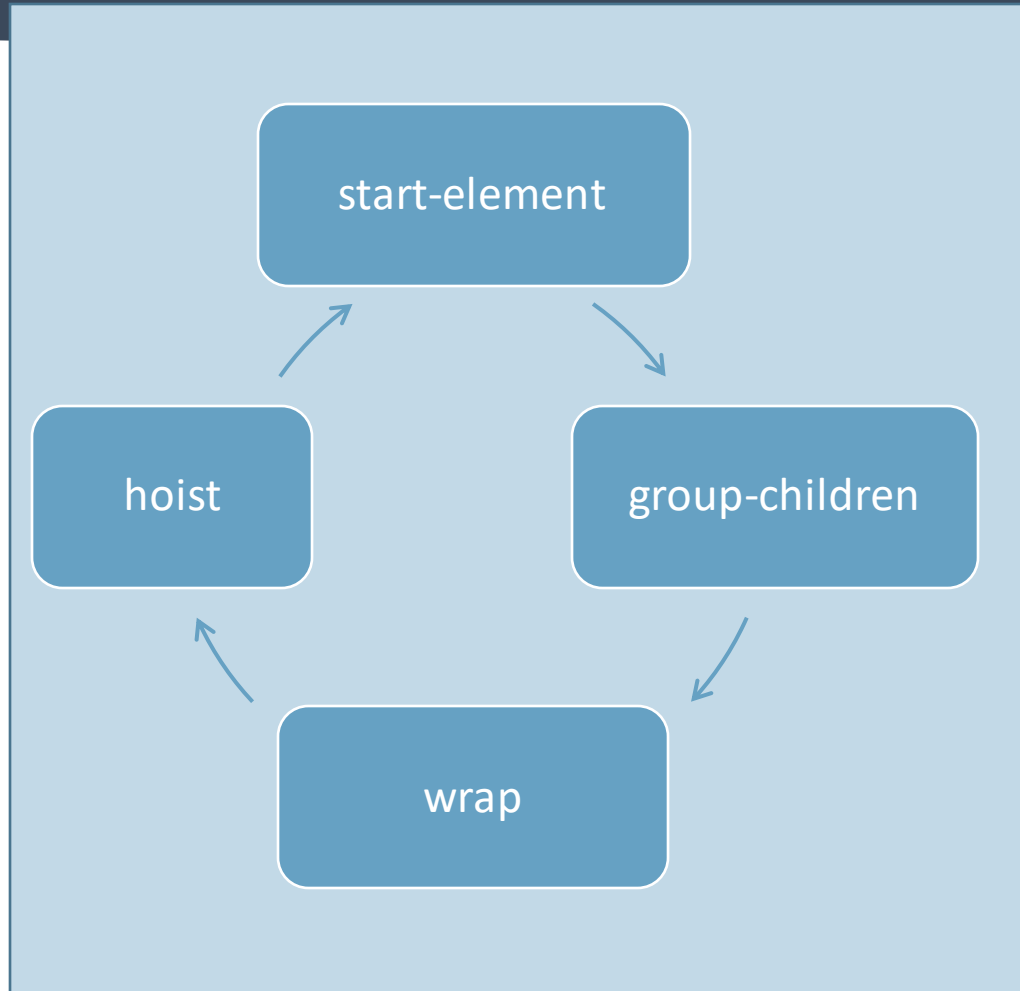
```
<xsl:variable name="expectedChildren"
  as="map(xs:string, xs:string+)"
  select="map {
    'table': (
      'caption','colgroup','tbody',
      'thead','tfoot'),
    'tbody': ('tr'),
    'thead': ('tr'),
    'tfoot': ('tr'),
    'tr': ('td','th'),
    'td': (''),
    'th': (''),
    'colgroup': ('col')
  }"/>
```

Context: Positioning the HTML Table Normalizer in the XSLT Pipeline

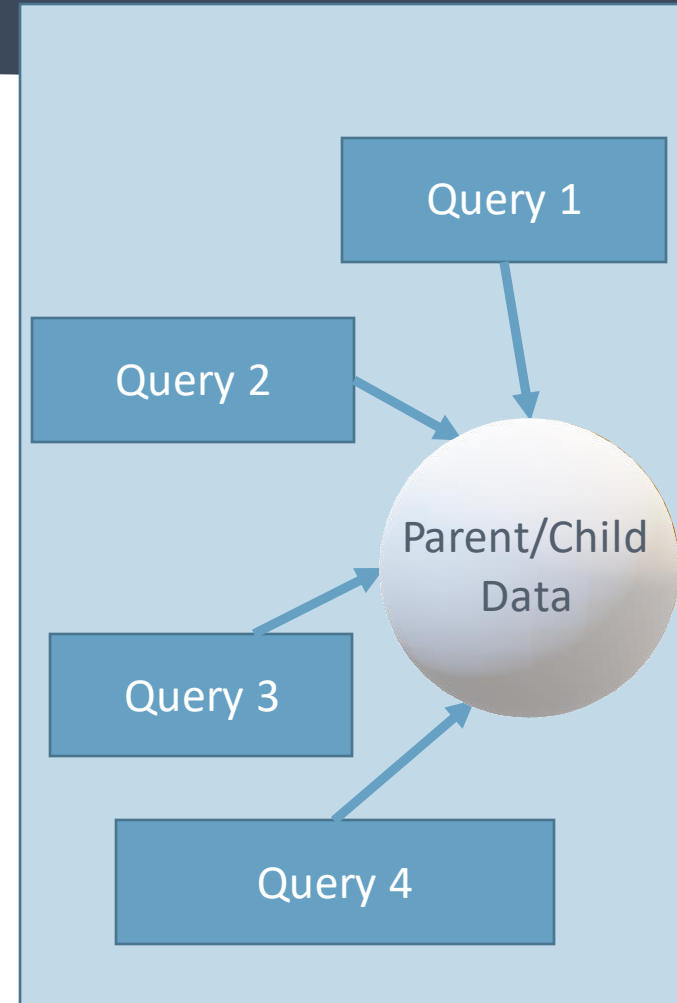


Normalizer: High-Level Design

Element-Tree Processing



Content Model

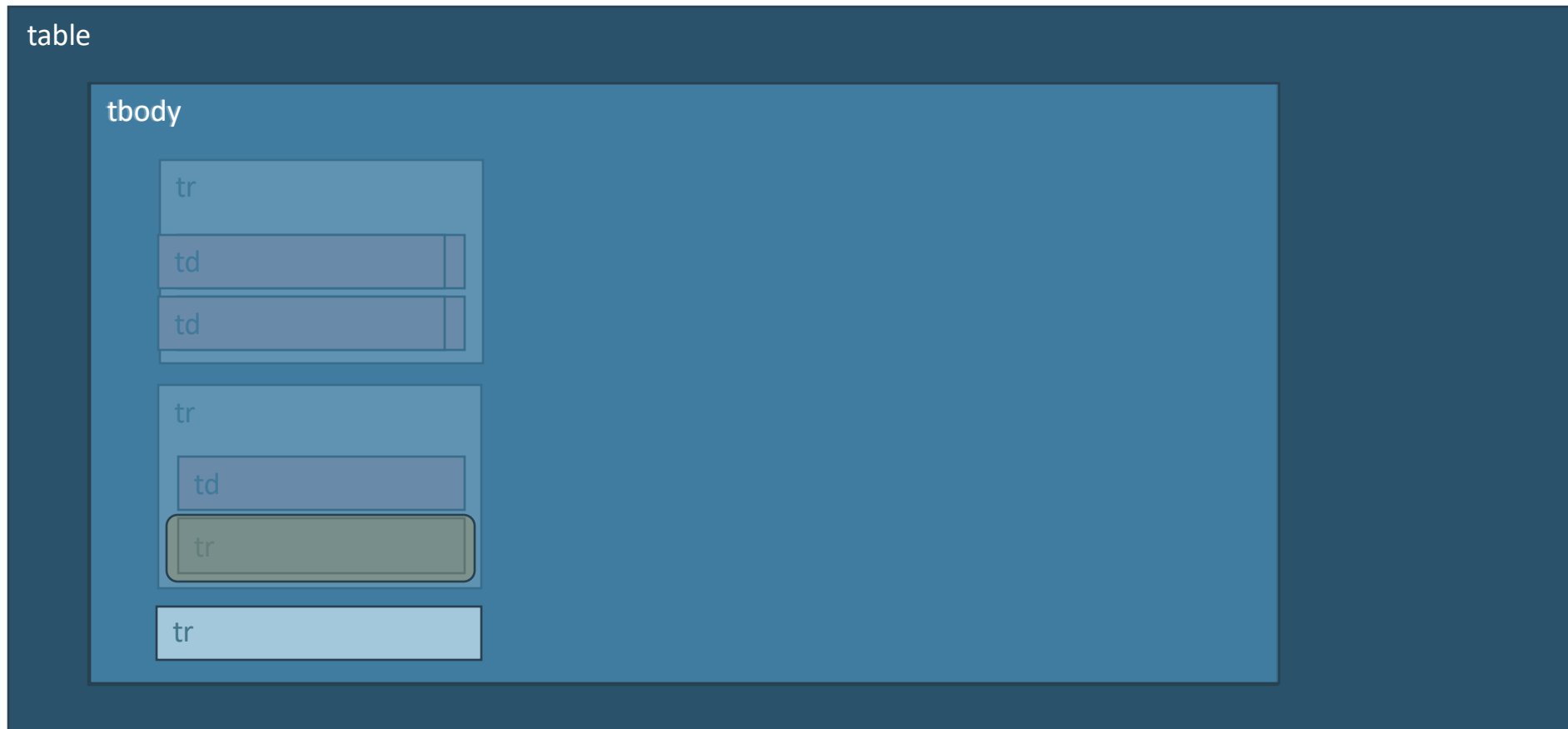


Element Wrapping



Content-Model Controlled Recursive Descent

10

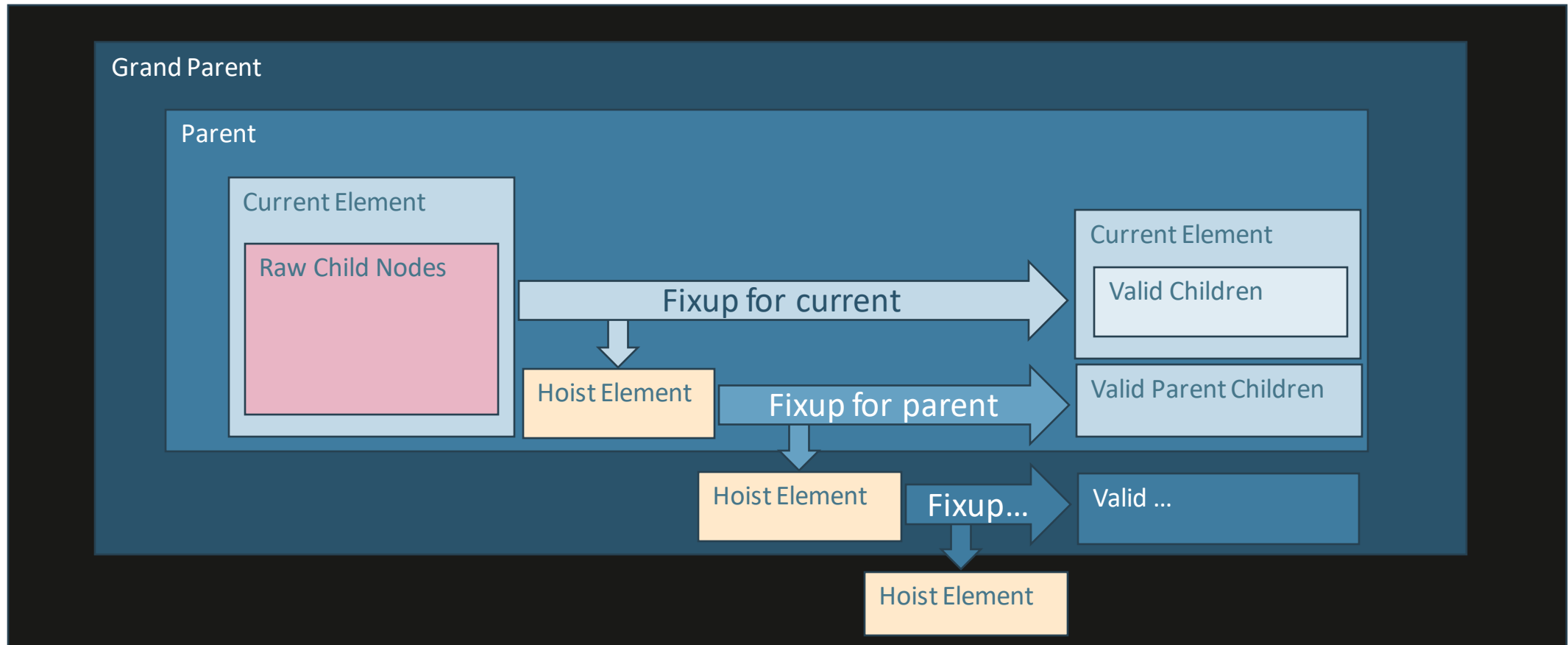


'Evading Capture' (Hoisting)



Recursive Descent with Element Hoisting

12



The xsl:template

‘Fix’ child nodes and insert as children and/or siblings of the current element

For each table element, assign wrapped child-nodes and hoist-position to variables

15

```
<xsl:template match="*[local-name()=$tableElementNames]" mode="inTable">
  <xsl:variable name="fixedChildNodes" as="node(*)" select="fn:fixupChildNodes(local-name(), node())"/>
  <xsl:variable name="firstHoistPosition" as="xs:integer" select="fn:firstHoistPosition($fixedChildNodes)"/>
  <xsl:choose>
    <xsl:when test="$firstHoistPosition > 0">
      <xsl:copy>
        <xsl:sequence select="@*, subsequence($fixedChildNodes, 1, $firstHoistPosition - 1)"/>
      </xsl:copy>
      <xsl:variable name="parentName" as="xs:string" select="fn:findParentName(local-name())"/>
      <xsl:variable name="hoistElementsForParent" as="node(*)">
        <xsl:apply-templates select="subsequence($fixedChildNodes, $firstHoistPosition)" mode="stripHoistWrapper"/>
      </xsl:variable>
      <xsl:message use-when="$debug" select="fn:msg-on-hoist(., $fixedChildNodes, $firstHoistPosition, $parentName,
        <xsl:sequence select="fn:fixupChildNodes($parentName, $hoistElementsForParent)"/>
      </xsl:when>
    <xsl:otherwise>
      <xsl:copy>
        <xsl:sequence select="@*, $fixedChildNodes"/>
      </xsl:copy>
    </xsl:otherwise>
  </xsl:choose>
</xsl:template>
```

Insert wrapped nodes before hoist-position into context-element copy

16

```
<xsl:template match="*[local-name() = $tableElementNames]" mode="inTable">
  <xsl:variable name="fixedChildNodes" as="node(*)" select="fn:fixupChildNodes(local-name(), node())"/>
  <xsl:variable name="firstHoistPosition" as="xs:integer" select="fn:firstHoistPosition($fixedChildNodes)"/>
  ..
  <xsl:choose>
    <xsl:when test="$firstHoistPosition > 0">
      <xsl:copy>
        <xsl:sequence select="@*, subsequence($fixedChildNodes, 1, $firstHoistPosition - 1)"/>
        </xsl:copy>
        <xsl:variable name="parentName" as="xs:string" select="fn:findParentName(local-name())"/>
        <xsl:variable name="hoistElementsForParent" as="node(*)">
          <xsl:apply-templates select="subsequence($fixedChildNodes, $firstHoistPosition)" mode="stripHoistWrapper"/>
        </xsl:variable>
        <xsl:message use-when="$debug" select="fn:msg-on-hoist(., $fixedChildNodes, $firstHoistPosition, $parentName,
        <xsl:sequence select="fn:fixupChildNodes($parentName, $hoistElementsForParent)"/>
        </xsl:when>
    <xsl:otherwise>
      <xsl:copy>
        <xsl:sequence select="@*, $fixedChildNodes"/>
      </xsl:copy>
    </xsl:otherwise>
  </xsl:choose>
</xsl:template>
```

Fix remaining nodes for the **context-parent** and append to result sequence

17

```
<xsl:template match="*[local-name() = $tableElementNames]" mode="inTable">
  <xsl:variable name="fixedChildNodes" as="node(*)" select="fn:fixupChildNodes(local-name(), node())"/>
  <xsl:variable name="firstHoistPosition" as="xs:integer" select="fn:firstHoistPosition($fixedChildNodes)"/>
  ..
  <xsl:choose>
    <xsl:when test="$firstHoistPosition > 0">
      <xsl:copy>
        <xsl:sequence select="@*, subsequence($fixedChildNodes, 1, $firstHoistPosition - 1)"/>
      </xsl:copy>
      <xsl:variable name="parentName" as="xs:string" select="fn:findParentName(local-name())"/>
      <xsl:variable name="hoistElementsForParent" as="node(*)">
        <xsl:apply-templates select="subsequence($fixedChildNodes, $firstHoistPosition)" mode="stripHoistWrapper"/>
      </xsl:variable>
      <xsl:message use-when="$debug" select="fn:msg-on-hoist(., $fixedChildNodes, $firstHoistPosition, $parentName,
        <xsl:sequence select="fn:fixupChildNodes($parentName, $hoistElementsForParent)"/>
      </xsl:when>
    <xsl:otherwise>
      <xsl:copy>
        <xsl:sequence select="@*, $fixedChildNodes"/>
      </xsl:copy>
    </xsl:otherwise>
  </xsl:choose>
</xsl:template>
```


The Node Wrapping Function

Group nodes for wrapping and invoke `<xsl:apply-templates>`

On the importance of range types and wrapping text in the different wrap-state

19

```
<xsl:function name="fn:fixupChildNodes" as="node()*">
  <xsl:param name="parentName" as="xs:string"/>
  <xsl:param name="childNodes" as="node()*"/>
  <xsl:iterate select="$childNodes">
    <xsl:param name="wrapperName" as="xs:string" select="''"/>
    <xsl:param name="wrapStart" as="xs:integer" select="0"/>

    <xsl:on-completion>
      <xsl:sequence select="fn:addWrapper($childNodes, $wrapperName, $wrapStart, count($childNodes) + 1)"/>
    </xsl:on-completion>

    <xsl:variable name="isWrapStarted" as="xs:boolean" select="$wrapStart gt 0"/>
    <xsl:variable name="nodeName" as="xs:string" select="local-name()"/>
    <xsl:variable name="validName" as="xs:string?" select="fn:findValidChild($nodeName, $parentName)"/>
    <xsl:variable name="state" as="map(xs:string, xs:boolean)"
      select="fn:calcWrapState(., $nodeName, $validName, $wrapperName, $isWrapStarted)"/>
    <xsl:if test="$state?hoist or $state?startWrap or $state?endWrap">
      <xsl:sequence select="fn:addWrapper($childNodes, $wrapperName, $wrapStart, position())"/>
    </xsl:if>
  </xsl:iterate>
</xsl:function>
```

...

On 'hoist', enclose all remaining nodes in 'hoist' element and break-iteration

20

```
<xsl:choose>
  ..<xsl:when test="$state?hoist">
    <deltaxml:hoist-nodes>
      ➡ <xsl:sequence select="subsequence($childNodes, position())"/>
    </deltaxml:hoist-nodes>
    <xsl:break/>
  </xsl:when>
  <xsl:when test="$state?startWrap">
    <xsl:next-iteration>
      ..<xsl:with-param name="wrapperName" select="$validName"/>
      ... <xsl:with-param name="wrapStart" select="position()"/>
    </xsl:next-iteration>
  </xsl:when>
  ..<xsl:when test="$state?endWrap">
    <xsl:apply-templates select="." mode="inTable"/>
    <xsl:next-iteration>
      ....<xsl:with-param name="wrapperName" select="''"/>
      <xsl:with-param name="wrapStart" select="0"/>
    </xsl:next-iteration>
  </xsl:when>

```

...

On 'startWrap' set the wrap-state: \$wrapperName and \$wrapStart

21

```
<xsl:choose>
  ..<xsl:when test="$state?hoist">
    <deltaxml:hoist-nodes>
      ....<xsl:sequence select="subsequence($childNodes, position())"/>
    </deltaxml:hoist-nodes>
    <xsl:break/>
  </xsl:when>
  <xsl:when test="$state?startWrap">
    <xsl:next-iteration>
      ..<xsl:with-param name="wrapperName" select="$validName"/>
      ..<xsl:with-param name="wrapStart" select="position()"/>
    </xsl:next-iteration>
  </xsl:when>
  ..<xsl:when test="$state?endWrap">
    <xsl:apply-templates select="." mode="inTable"/>
    <xsl:next-iteration>
      ....<xsl:with-param name="wrapperName" select="''"/>
      <xsl:with-param name="wrapStart" select="0"/>
    </xsl:next-iteration>
  </xsl:when>

```

...

On 'endWrap', apply templates to current node and reset the wrap-state

22

```
<xsl:choose>
  ..<xsl:when test="$state?hoist">
    <deltaxml:hoist-nodes>
      ....<xsl:sequence select="subsequence($childNodes, position())"/>
    </deltaxml:hoist-nodes>
    <xsl:break/>
  </xsl:when>
  <xsl:when test="$state?startWrap">
    <xsl:next-iteration>
      ..<xsl:with-param name="wrapperName" select="$validName"/>
      ..<xsl:with-param name="wrapStart" select="position()"/>
    </xsl:next-iteration>
  </xsl:when>
  ..<xsl:when test="$state?endWrap">
    ➡ <xsl:apply-templates select="." mode="inTable"/>
    <xsl:next-iteration>
      ➡ <xsl:with-param name="wrapperName" select="''"/>
      <xsl:with-param name="wrapStart" select="0"/>
    </xsl:next-iteration>
  </xsl:when>
</xsl:choose>
```

...

Otherwise (state has not changed),
<xsl:apply-templates> to current-node, if not 'buffering'

23

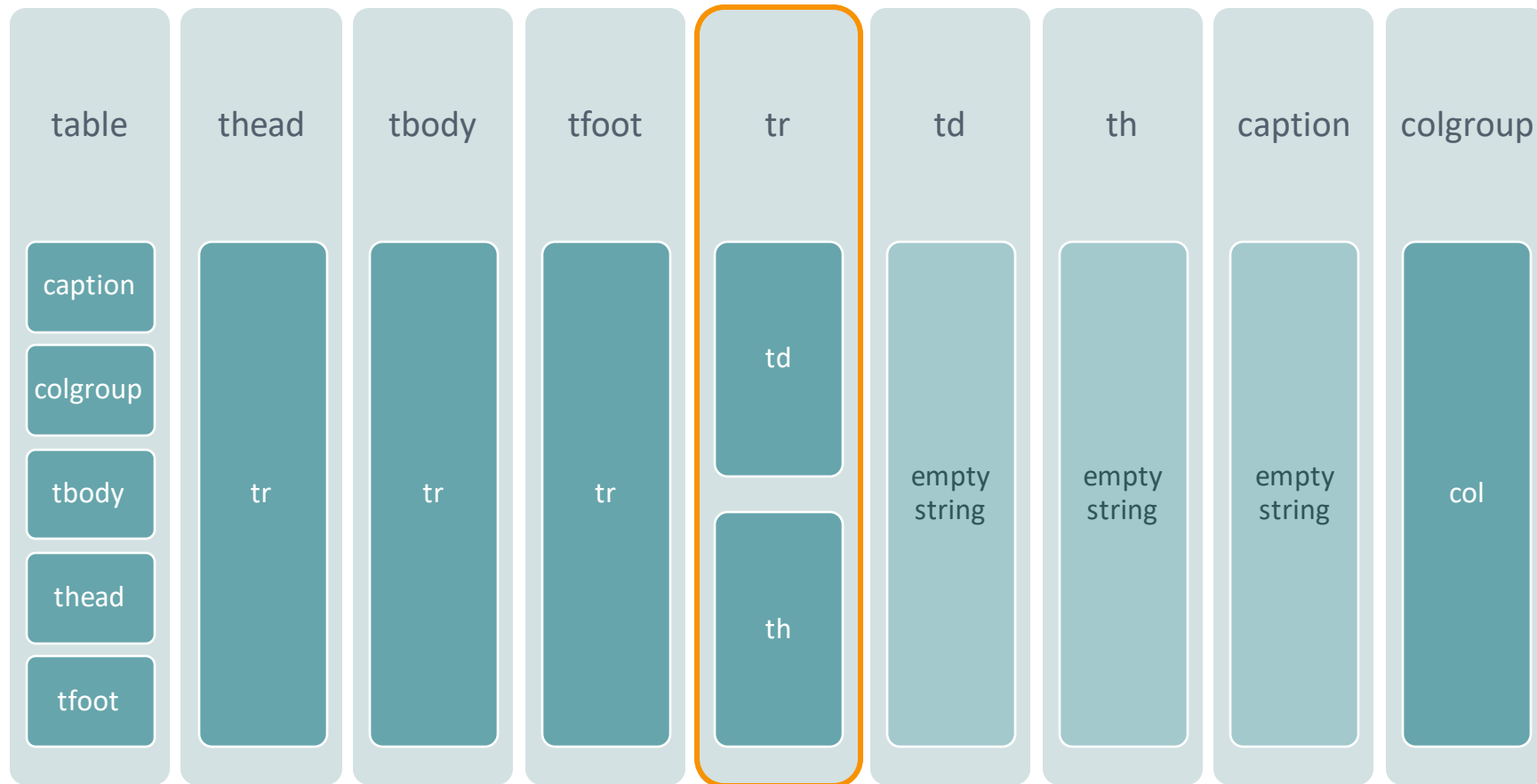
...

```
... <xsl:otherwise>
...   <xsl:apply-templates select="if ($isWrapStarted) then () else ." mode="inTable"/>
... </xsl:otherwise>
</xsl:choose>
```

The Content Model

The 'Parent / Child ' Data Map

`<xsl:variable name="expectedChildren" as="map(xs:string, xs:string+)"...`

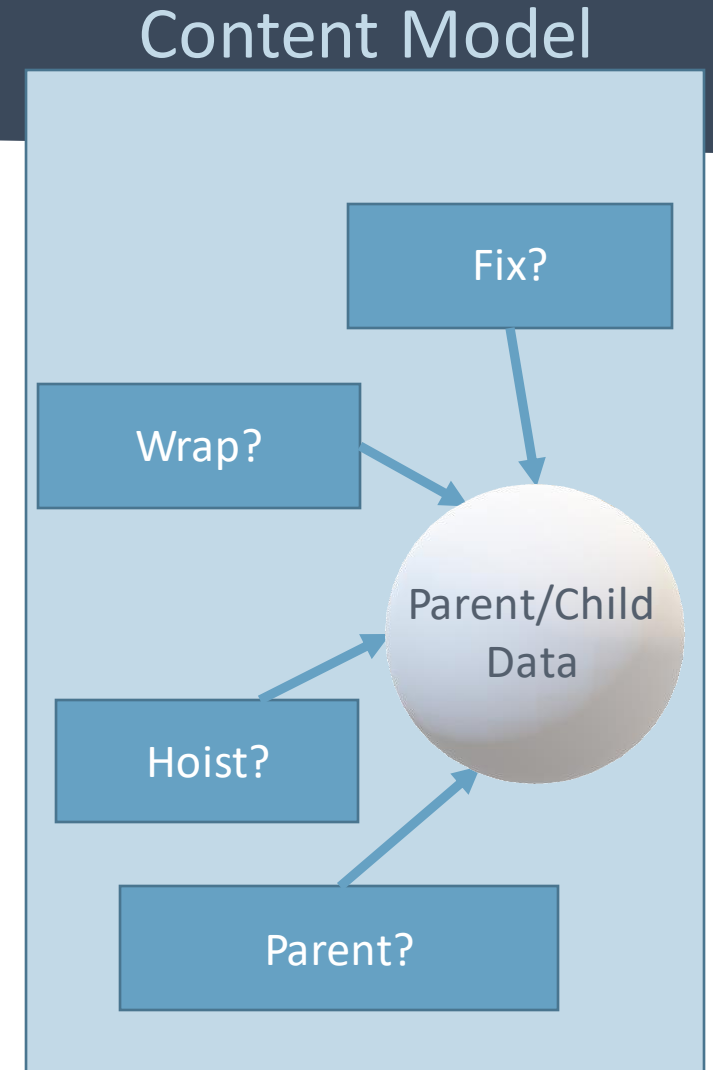


Content Model Function Summary

26

For any XML element:

- **Fix?** Is it in the table-elements list?
- **Wrap?** What wrapper is expected, given the parent element?
- **Hoist?** Must this be hoisted, given the parent element?
- **Parent?** What is the valid parent element for this table-element?



Tracing XSLT Execution

Using `<xsl:message>` with coloring and formatting

#2 Wrap `<table>` in `<tbody>`

28

----- Insertion Mode: `table` -----

1 of 5 nodes of <code>table</code>	name: <code>tr</code>	expected: <code>tbody</code>	buffer: start	<code><tbody></code>
2 of 5 nodes of <code>table</code>	name: <code>td</code>	expected: <code>tbody</code>	buffer: continue	<code>...</code>
3 of 5 nodes of <code>table</code>	name: <code>td</code>	expected: <code>tbody</code>	buffer: continue	<code>...</code>
4 of 5 nodes of <code>table</code>	name: <code>td</code>	expected: <code>tbody</code>	buffer: continue	<code>...</code>
5 of 5 nodes of <code>table</code>	name: <code>tr</code>	expected: <code>tbody</code>	buffer: continue	<code>...</code>
completed: (wrap buffer)				<code></tbody></code>

----- Insertion Mode: `tbody` -----

1 of 5 nodes of <code>tbody</code>	name: <code>tr</code>	expected: <code>tr</code>	buffer: no	<code>...</code>
------------------------------------	-----------------------	---------------------------	------------	------------------

----- Insertion Mode: `tr` -----

1 of 3 nodes of <code>tr</code>	name: <code>td</code>	expected: <code>td</code>	buffer: no	<code>..</code>
2 of 3 nodes of <code>tr</code>	name: <code>td</code>	expected: <code>td</code>	buffer: no	<code>..</code>
3 of 3 nodes of <code>tr</code>	name: <code>td</code>	expected: <code>td</code>	buffer: no	<code>..</code>
completed				

2 of 5 nodes of <code>tbody</code>	name: <code>td</code>	expected: <code>tr</code>	buffer: start	<code><tr></code>
3 of 5 nodes of <code>tbody</code>	name: <code>td</code>	expected: <code>tr</code>	buffer: continue	<code>...</code>
4 of 5 nodes of <code>tbody</code>	name: <code>td</code>	expected: <code>tr</code>	buffer: continue	<code>...</code>
5 of 5 nodes of <code>tbody</code>	name: <code>tr</code>	expected: <code>tr</code>	buffer: end	<code></tr></code>

----- Insertion Mode: `tr` -----

1 of 3 nodes of <code>tr</code>	name: <code>td</code>	expected: <code>td</code>	buffer: no	<code>..</code>
2 of 3 nodes of <code>tr</code>	name: <code>td</code>	expected: <code>td</code>	buffer: no	<code>..</code>
3 of 3 nodes of <code>tr</code>	name: <code>td</code>	expected: <code>td</code>	buffer: no	<code>..</code>
completed				

----- Insertion Mode: `tr` -----

1 of 3 nodes of <code>tr</code>	name: <code>td</code>	expected: <code>td</code>	buffer: no	<code>...</code>
2 of 3 nodes of <code>tr</code>	name: <code>td</code>	expected: <code>td</code>	buffer: no	<code>...</code>
3 of 3 nodes of <code>tr</code>	name: <code>td</code>	expected: <code>td</code>	buffer: no	<code>...</code>

Input XML

```
<html lang="en">
<title>test</title>
<table>
  <tr>
    <td>half</td>
    <td>quarter</td>
    <td>third</td>
  </tr>
  <td>one</td>
  <td>two</td>
  <td>three</td>
  <tr>
    <td>ten</td>
    <td>twenty</td>
    <td>thirty</td>
  </tr>
</table>
</html>
```

Simple Element Hoisting Example

```

----- Insertion Mode: table -----
1 of 1 nodes of table  name: tbody  expected: tbody  buffer: no  ...

----- Insertion Mode: tbody -----
1 of 3 nodes of tbody  name: tr  expected: tr  buffer: no  ...

----- Insertion Mode: tr -----
1 of 1 nodes of tr  name: td  expected: td  buffer: no  ...
completed
2 of 3 nodes of tbody  name: tbody  expected:  buffer: hoist  <hoist/>
== break iteration in mode tbody ==
Insert in tbody:  alpha  <tr />
Hoist to table:  (
                 bravo  <tbody />,
                 charlie <tr />)

----- Insertion Mode: table -----
1 of 2 nodes of table  name: tbody  expected: tbody  buffer: no  ...

----- Insertion Mode: tbody -----
1 of 1 nodes of tbody  name: tr  expected: tr  buffer: no  ...
2 of 2 nodes of table  name: tr  expected: tbody  buffer: start  <tbody>
completed: (wrap buffer)  </tbody>

```

```

<html>
  <table>
    <tbody>
      <tr>
        <td>alpha</td>
      </tr>
      <tbody>
        <tr>
          <td>bravo</td>
        </tr>
      </tbody>
      <tr>
        <td>charlie</td>
      </tr>
    </tbody>
  </table>
</html>

```

Testing XSLT with XSpec (in VS Code)

```

PROBLEMS  OUTPUT  DEBUG CONSOLE  DOCUMENTATION  TIMELINE  TERMINAL

XSpec testing with SAXON HE 11.3

missing wrapper elements
✓ valid tbody
✓ missing tbody
✓ single row table missing tbody and tr

hoist th element from td element
✗ hoisted th element and trailing th elements require tr wrapper (snip)
-----

hoist tbody from td
✓ with td inside thead
✓ multi-hoist: tbody within tr within td

hoist caption from tbody element
✗ caption between tr element

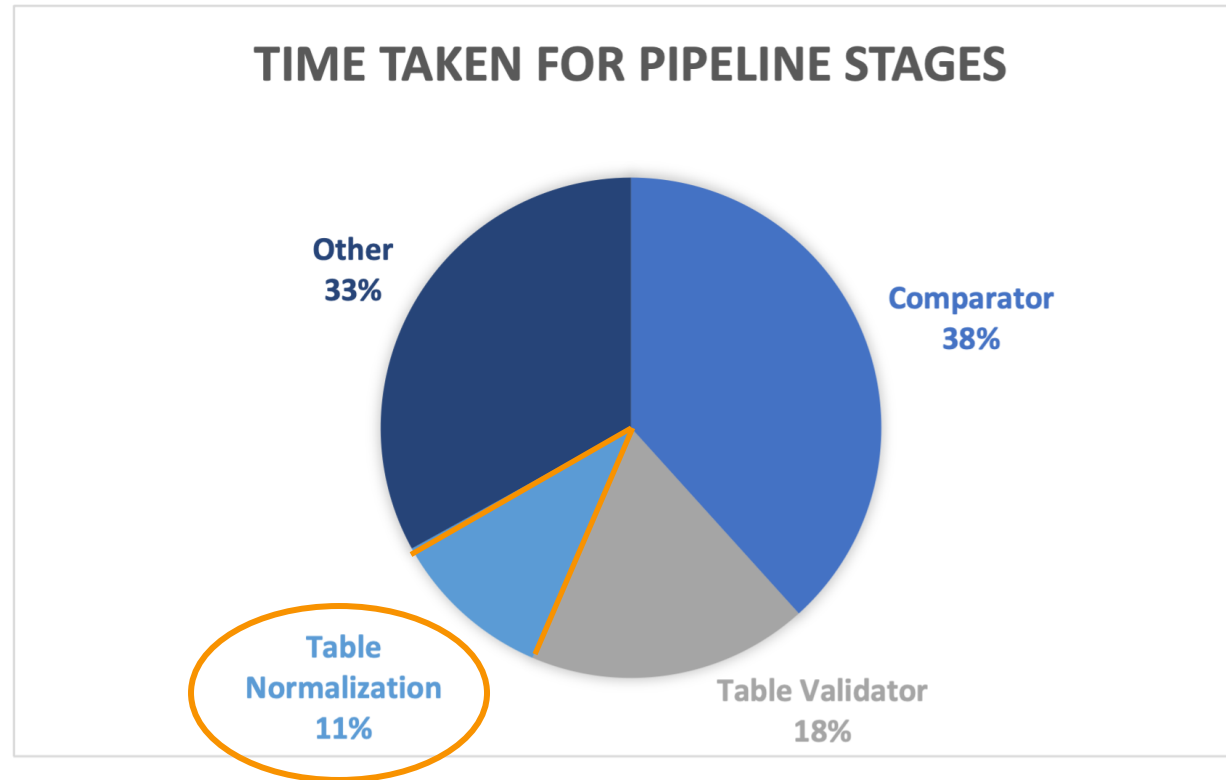
-----
passed:  24
pending: 0
failed:  2
-----
total:   26

```

Analysing Failed XSpec Tests

Expected		Actual	
416-	<td/>	→ 416+	<td>
417-	</tr>	417+	<tr>
418-	<tr>	418+	<th>c1</th>
419-	<th>c1</th>	419+	<th>c2</th>
420-	<th>c2</th>	420+	<th>c3</th>
421-	<th>c3</th>	421+	<th>c4</th>
422-	<th>c4</th>	422+	</tr>
423-	</tr>	423+	<deltaxml:hoist-nodes>
424-	</thead>	424+	<tbody>
425-	<tbody>	425+	<td>new</td>
426-	<tr>	426+	</tbody>
427-	<td>new</td>	427+	</deltaxml:hoist-nodes>
428-	</tr>	428+	<th>No.</th>
429-	</tbody>	429+	<th>Name</th>
430-	<tbody>	430+	<th>Email</th>
431-	<tr>	431+	<th>DOB</th>
432-	<th>No.</th>	432+	</td>
433-	<th>Name</th>		
434-	<th>Email</th>		
435-	<th>DOB</th>		

Performance



Summary

- Benefits of a Strict Content Model within a Pipeline
- Separation of Content Model from XML Structure Processing
- Recursive element 'Wrapping' during descent
- 'Hoisting' achieved by 'Evading Capture' during descent
- Tracing recursive execution with `<xsl:message/>`

<xsl:on-completion/>