

# The Streaming Complexity of Validating XML Documents

Christian Konrad

`christian.konrad@liafa.jussieu.fr`

and

Frédéric Magniez

`frederic.magniez@liafa.jussieu.fr`



**LIAFA**

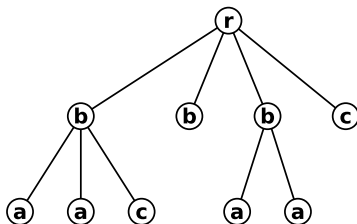
University Paris Diderot - Paris 7  
Paris

**ICDT 2012**

# What is XML?

- **XML document:** sequence of opening and closing tags

```
<r>
  <b>
    <a</a>
    <a</a>
    <c</c>
  </b>
  <b</b>
  <b>
    <a</a>
    <a</a>
  </b>
  <c</c>
</r>
```



**Notation:**  $rb\bar{a}\bar{a}\bar{a}\bar{a}\bar{c}\bar{c}\bar{b}\bar{b}\bar{b}\bar{b}\bar{a}\bar{a}\bar{a}\bar{a}\bar{b}\bar{c}\bar{c}\bar{r}$

$\text{pos}(a), \text{pos}(\bar{a})$ : position in XML document

$\text{depth}(a), \text{depth}(\bar{a})$ : depth of corresp. node

- **Depth first tree traversal:** down step gives opening tag, up step gives closing tag

**Well-formedness:** An XML document is well-formed iff each opening tag is closed by its corresponding closing tag

- $ra\bar{a}b\bar{b}\bar{r}$  is well-formed
- $ra\bar{b}b\bar{a}\bar{r}$  is not well-formed

Only well-formed documents correspond to a tree

# Well-formedness and Validity

**Well-formedness:** An XML document is well-formed iff each opening tag is closed by its corresponding closing tag

- $ra\bar{a}b\bar{b}\bar{r}$  is well-formed
- $rab\bar{b}\bar{a}\bar{r}$  is not well-formed

Only well-formed documents correspond to a tree

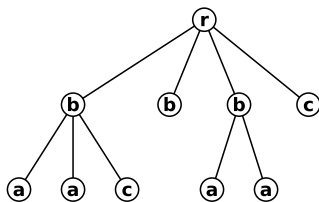
**Validity:** is checked wrt. a DTD (**D**ocument **T**ype **D**efinition)

$$r \rightarrow b^*c^+$$

$$b \rightarrow a^*c?|\epsilon$$

$$a \rightarrow \epsilon$$

$$c \rightarrow \epsilon$$



**Difficulty:** relate each label to labels of its children

# Well-formedness and Validity

**Well-formedness:** An XML document is well-formed iff each opening tag is closed by its corresponding closing tag

- $ra\bar{a}b\bar{b}\bar{r}$  is well-formed
- $rab\bar{b}\bar{a}\bar{r}$  is not well-formed

Only well-formed documents correspond to a tree

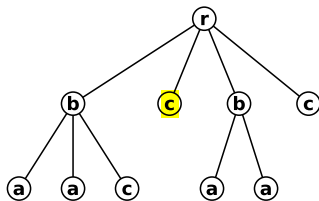
**Validity:** is checked wrt. a DTD (**D**ocument **T**ype **D**efinition)

$$r \rightarrow b^*c^+$$

$$b \rightarrow a^*c?|\epsilon$$

$$a \rightarrow \epsilon$$

$$c \rightarrow \epsilon$$



**Difficulty:** relate each label to labels of its children

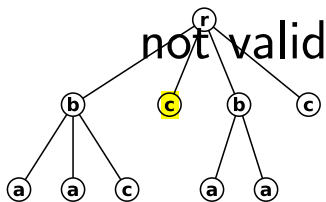
# Well-formedness and Validity

**Well-formedness:** An XML document is well-formed iff each opening tag is closed by its corresponding closing tag

- $ra\bar{a}b\bar{b}\bar{r}$  is well-formed
- $rab\bar{b}\bar{a}\bar{r}$  is not well-formed

Only well-formed documents correspond to a tree

**Validity:** is checked wrt. a DTD (**D**ocument **T**ype **D**efinition)

$$\begin{aligned} r &\rightarrow b^*c^+ \\ b &\rightarrow a^*c?|\epsilon \\ a &\rightarrow \epsilon \\ c &\rightarrow \epsilon \end{aligned}$$


**Difficulty:** relate each label to labels of its children

# Stream Computation

- **Objective:** compute some function  $f(x_1, \dots, x_n)$  given only sequential access



$x_1$   $x_2$   $x_3$   $x_4$   $x_5$   $x_6$  ...  $x_n$

How much RAM is required for the computation of  $f$ ?

# Stream Computation

- **Objective:** compute some function  $f(x_1, \dots, x_n)$  given only sequential access



$x_1$   $x_2$   $x_3$   $x_4$   $x_5$   $x_6$   $\dots$   $x_n$


How much RAM is required for the computation of  $f$ ?



# Stream Computation

- **Objective:** compute some function  $f(x_1, \dots, x_n)$  given only sequential access

$x_1$   $x_2$   $x_3$   $x_4$   $x_5$   $x_6$   $\dots$   $x_n$



How much RAM is required for the computation of  $f$ ?

# Stream Computation

- **Objective:** compute some function  $f(x_1, \dots, x_n)$  given only sequential access

$x_1$   $x_2$   $x_3$   $x_4$   $x_5$   $x_6$   $\dots$   $x_n$



How much RAM is required for the computation of  $f$ ?

- **Motivation:** massive data sets
  - Storage on external disks, cheap sequential access
  - Data streams over the internet
  - XML databases can be huge

# Stream Computation

- **Objective:** compute some function  $f(x_1, \dots, x_n)$  given only sequential access

$x_1$   $x_2$   $x_3$   $x_4$   $x_5$   $x_6$  ...  $x_n$



How much RAM is required for the computation of  $f$ ?

- **Motivation:** massive data sets
  - Storage on external disks, cheap sequential access
  - Data streams over the internet
  - XML databases can be huge
- **Scenarios:**
  - multiple passes
  - deterministic/randomized
  - bidirectional
  - ...
  - **auxiliary streams** (external memory)

# Auxiliary Streams (external memory)

- **Example:** Merge Sort with 3 streams,  $O(\log N)$  passes,  $O(\log N)$  space

Stream 1:  $x_1$   $x_2$   $x_3$  ...  $x_n$

Stream 2:

Stream 3:

input on stream 1

## Auxiliary Streams (external memory)

- **Example:** Merge Sort with 3 streams,  $O(\log N)$  passes,  $O(\log N)$  space

Stream 1:  $x_1$   $x_2$   $x_3$  ...  $x_n$   
Stream 2:  $x_1$   $x_3$  ...  $x_{n-1}$   
Stream 3:  $x_2$   $x_4$  ...  $x_n$

copy numbers alternately onto stream 2 and stream 3

## Auxiliary Streams (external memory)

- **Example:** Merge Sort with 3 streams,  $O(\log N)$  passes,  $O(\log N)$  space

Stream 1:	$x_1$	$x_2$	$x_3$	$\dots$	$x_n$
Stream 2:	$X_1$	$X_3$	$\dots$	$X_{n-1}$	
Stream 3:	$X_2$	$X_4$	$\dots$	$X_n$	

think of numbers as sorted blocks of size 1

## Auxiliary Streams (external memory)

- **Example:** Merge Sort with 3 streams,  $O(\log N)$  passes,  $O(\log N)$  space

Stream 1:  $X_{12}$   $X_{34}$   $\dots$   $X_{n-1,n}$   
Stream 2:  $X_1$   $X_3$   $\dots$   $X_{n-1}$   
Stream 3:  $X_2$   $X_4$   $\dots$   $X_n$

*merge operation:* merge blocks into blocks of size 2 onto stream 1

## Auxiliary Streams (external memory)

- **Example:** Merge Sort with 3 streams,  $O(\log N)$  passes,  $O(\log N)$  space

Stream 1:  $X_{12}$   $X_{34}$  ...  $X_{n-1,n}$   
Stream 2:  $X_{12}$   $X_{56}$  ...  
Stream 3:  $X_{34}$   $X_{78}$  ...

copy blocks of size 2 alternately onto stream 2 and stream 3



## Auxiliary Streams (external memory)

- **Example:** Merge Sort with 3 streams,  $O(\log N)$  passes,  $O(\log N)$  space

Stream 1:  $X_{1234}$   $X_{5678}$  ...  $X_{n-3,\dots,n}$   
Stream 2:  $X_{12}$   $X_{56}$  ...  
Stream 3:  $X_{34}$   $X_{78}$  ...

*merge operation:* merge blocks of size 2 into blocks of size 4 onto stream 1

## Auxiliary Streams (external memory)

- **Example:** Merge Sort with 3 streams,  $O(\log N)$  passes,  $O(\log N)$  space

Stream 1: ...

Stream 2: ...

Stream 3: ...

repeat this procedure until we obtain a sorted block of size  $n$

## Auxiliary Streams (external memory)

- **Example:** Merge Sort with 3 streams,  $O(\log N)$  passes,  $O(\log N)$  space

Stream 1:  $X_{1\dots n}$

Stream 2: ...

Stream 3: ...

repeat this procedure until we obtain a sorted block of size  $n$

# Auxiliary Streams (external memory)

- **Example:** Merge Sort with 3 streams,  $O(\log N)$  passes,  $O(\log N)$  space

Stream 1:  $X_{1\dots n}$

Stream 2: ...

Stream 3: ...

constant number of passes to double block size  $\rightarrow O(\log N)$  passes

# Auxiliary Streams (external memory)

- **Example:** Merge Sort with 3 streams,  $O(\log N)$  passes,  $O(\log N)$  space

Stream 1:  $X_{1\dots n}$

Stream 2: ...

Stream 3: ...

constant number of passes to double block size  $\rightarrow O(\log N)$  passes

- **Important parameters:**
  - $k(N)$  auxiliary streams  
usually in addition to one read-only input stream
  - $p(N)$  passes
  - $s(N)$  random access space

## Well-formedness: Reduction to DYCK languages

- **DYCK(k)**: well-parenthesized words,  $k$  types of parenthesis  
 $((())) \in \text{DYCK}(2)$ ,  $(\{\}\}) \in \text{DYCK}(3)$

# Well-formedness: Reduction to DYCK languages

- **DYCK(k)**: well-parenthesized words,  $k$  types of parenthesis

$([()]) \in \text{DYCK}(2)$ ,  $([\{\}]) \in \text{DYCK}(3)$

- **Well-formedness**: document well-formed if in  $\text{DYCK}(k)$ :

$rb\bar{a}\bar{a}\bar{a}\bar{c}\bar{c}\bar{b}\bar{b}\bar{b}\bar{a}\bar{a}\bar{a}\bar{b}\bar{c}\bar{c}\bar{r}$   
 $(r(b(a)a(a)a(c)c)b(b)b(b(a)a(a)a)b(c)c)r$

# Well-formedness: Reduction to DYCK languages

- **DYCK(k)**: well-parenthesized words,  $k$  types of parenthesis

$$([()]) \in \text{DYCK}(2), ([\{\}]) \in \text{DYCK}(3)$$

- **Well-formedness**: document well-formed if in  $\text{DYCK}(k)$ :

$$rba\bar{a}a\bar{a}c\bar{c}b\bar{b}b\bar{b}a\bar{a}a\bar{a}b\bar{c}c\bar{r}$$
$$(r(b(a)a(a)a(c)c)b(b)b(b(a)a(a)a)b(c)c)r$$

- **Streaming Algorithms**: Checking DYCK membership

Theorem (F. Magniez, C. Mathieu, and A. Nayak, STOC 2010)

- There is a **randomized 1-pass algorithm** that decides membership to  $\text{DYCK}(k)$  with space  $O(\sqrt{N \log k \log(N \log k)})$ .
- There is a **bidirectional randomized 2-passes algorithm** that decides membership to  $\text{DYCK}(k)$  with space  $O((\log(N \log k))^2)$ .



# Well-formedness: Reduction to DYCK languages

- **DYCK(k)**: well-parenthesized words,  $k$  types of parenthesis

$$([()]) \in \text{DYCK}(2), ([\{\}]) \in \text{DYCK}(3)$$

- **Well-formedness**: document well-formed if in  $\text{DYCK}(k)$ :

$$rba\bar{a}a\bar{a}c\bar{c}b\bar{b}b\bar{b}a\bar{a}a\bar{a}b\bar{c}c\bar{r}$$
$$(r(b(a)a(a)a(c)c)b(b)b(b(a)a(a)a)b(c)c)r$$

- **Streaming Algorithms**: Checking DYCK membership

Theorem (F. Magniez, C. Mathieu, and A. Nayak, STOC 2010)

- There is a **randomized 1-pass algorithm** that decides membership to  $\text{DYCK}(k)$  with space  $O(\sqrt{N \log k \log(N \log k)})$ .
- There is a **bidirectional randomized 2-passes algorithm** that decides membership to  $\text{DYCK}(k)$  with space  $O((\log(N \log k))^2)$ .

From now on: **XML documents are well-formed**

# Starting Point

- **Prior works:** [Segoufin Sirangelo, 07], [Segoufin Vianu, 02]  
Characterization of DTDs that allow deterministic constant space validation in 1-pass
- **Upper bound:** stack based algorithm, space linear to depth of document, 1-pass deterministic
- **Lower bound:** ternary trees: any  $p$  pass randomized streaming algorithm deciding validity requires  $\Omega(N/p)$  space

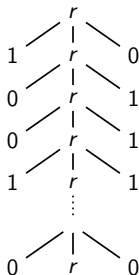
# Starting Point

- **Prior works:** [Segoufin Sirangelo, 07], [Segoufin Vianu, 02]  
Characterization of DTDs that allow deterministic constant space validation in 1-pass
- **Upper bound:** stack based algorithm, space linear to depth of document, 1-pass deterministic
- **Lower bound:** ternary trees: any  $p$  pass randomized streaming algorithm deciding validity requires  $\Omega(N/p)$  space

**DTD:**

$r \rightarrow 0r1 \mid 1r0 \mid 0r0 \mid \epsilon$

$0, 1 \rightarrow \epsilon$



$r1\bar{1}r0\bar{0}r0\bar{0}r1\bar{1}r \dots 0\bar{0}r\bar{r}0\bar{0} \dots \bar{r}1\bar{1}\bar{r}1\bar{1}\bar{r}1\bar{1}\bar{r}0\bar{0}\bar{r}$

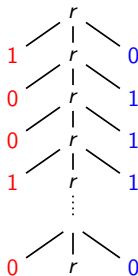
# Starting Point

- **Prior works:** [Segoufin Sirangelo, 07], [Segoufin Vianu, 02]  
Characterization of DTDs that allow deterministic constant space validation in 1-pass
- **Upper bound:** stack based algorithm, space linear to depth of document, 1-pass deterministic
- **Lower bound:** ternary trees: any  $p$  pass randomized streaming algorithm deciding validity requires  $\Omega(N/p)$  space

**DTD:**

$r \rightarrow 0r1 \mid 1r0 \mid 0r0 \mid \epsilon$

$0, 1 \rightarrow \epsilon$



$r1\bar{1}r0\bar{0}r0\bar{0}r1\bar{1}r \dots 0\bar{0}r\bar{r}0\bar{0} \dots \bar{r}1\bar{1}\bar{r}1\bar{1}\bar{r}1\bar{1}\bar{r}0\bar{0}\bar{r}$

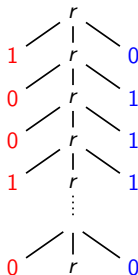
# Starting Point

- **Prior works:** [Segoufin Sirangelo, 07], [Segoufin Vianu, 02]  
Characterization of DTDs that allow deterministic constant space validation in 1-pass
- **Upper bound:** stack based algorithm, space linear to depth of document, 1-pass deterministic
- **Lower bound:** ternary trees: any  $p$  pass randomized streaming algorithm deciding validity requires  $\Omega(N/p)$  space

**DTD:**

$r \rightarrow 0r1 \mid 1r0 \mid 0r0 \mid \epsilon$

$0, 1 \rightarrow \epsilon$



$r1\bar{1}r0\bar{0}r0\bar{0}r1\bar{1}r \dots 0\bar{0}r\bar{r}0\bar{0} \dots \bar{r}1\bar{1}\bar{r}1\bar{1}\bar{r}1\bar{1}\bar{r}0\bar{0}\bar{r}$

- **Reduction:** Set-Disjointness in Communication Complexity

## Theorem

*There is a bidirectional  $O(\log N)$ -pass **deterministic** streaming algorithm for validity of arbitrary XML files and arbitrary DTDs with space  $O(\log^2 N)$  and 3 auxiliary streams.*

## Theorem

*There is a bidirectional  $O(\log N)$ -pass **deterministic** streaming algorithm for validity of arbitrary XML files and arbitrary DTDs with space  $O(\log^2 N)$  and 3 auxiliary streams.*

### Steps:

- 1 **Using 3 aux. streams,  $O(\log N)$  space,  $O(\log N)$  passes:**  
Compute the FCNS (First-Child-Next-Sibling) encoding of the original document (encoding as a binary tree)
- 2 **Using 2 bidirectional passes,  $O(\log^2 N)$  space:**  
Check validity based on this binary encoding  
Algorithm inspired by algorithm for checking validity of binary trees

## Theorem

*There is a one-pass **deterministic** algorithm using  $O(\sqrt{N \log N})$  space for checking validity of binary trees.*

**Conjecture:** there is no one-pass algorithm using  $o(\sqrt{N \log N})$  space even when randomization is allowed

## Theorem

*There is a bidirectional two-passes **deterministic** algorithm using  $O(\log^2 N)$  space for checking validity of binary trees.*



# Two-passes Algorithm for Validity of binary Trees

## Lemma (1)

*There is a one-pass **deterministic** algorithm using  $O(\log^2 N)$  space that verifies validity of all nodes which have a left subtree that is at least as large as its right subtree.*

# Two-passes Algorithm for Validity of binary Trees

## Lemma (1)

*There is a one-pass **deterministic** algorithm using  $O(\log^2 N)$  space that verifies validity of all nodes which have a left subtree that is at least as large as its right subtree.*

## Theorem

*There is a bidirectional two-passes **deterministic** algorithm using  $O(\log^2 N)$  space for validity.*

# Two-passes Algorithm for Validity of binary Trees

## Lemma (1)

*There is a one-pass **deterministic** algorithm using  $O(\log^2 N)$  space that verifies validity of all nodes which have a left subtree that is at least as large as its right subtree.*

## Theorem

*There is a bidirectional two-passes **deterministic** algorithm using  $O(\log^2 N)$  space for validity.*

### Proof:

- 1 Run algorithm of Lemma 1

# Two-passes Algorithm for Validity of binary Trees

## Lemma (1)

*There is a one-pass **deterministic** algorithm using  $O(\log^2 N)$  space that verifies validity of all nodes which have a left subtree that is at least as large as its right subtree.*

## Theorem

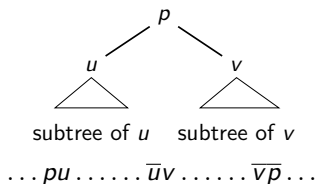
*There is a bidirectional two-passes **deterministic** algorithm using  $O(\log^2 N)$  space for validity.*

### Proof:

- 1 Run algorithm of Lemma 1
- 2 Run algorithm of Lemma 1 on the input stream read from right to left interpreting opening tags as closing tags and vice versa.

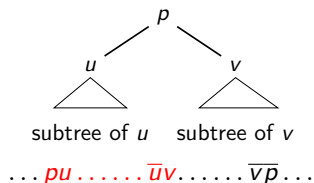
# Binary Trees - general ideas

- **Goal:** for all internal nodes  $p$ : relate  $p$  to its children  $u, v$  via  $\text{check}(p, u, v)$



# Binary Trees - general ideas

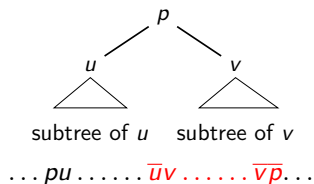
- **Goal:** for all internal nodes  $p$ : relate  $p$  to its children  $u, v$  via  $\text{check}(p, u, v)$



- Two chances for verification:
  - Top down using  $p, \bar{u}, v$

# Binary Trees - general ideas

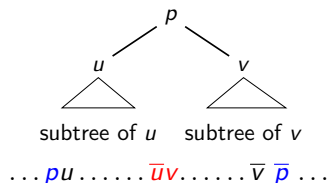
- **Goal:** for all internal nodes  $p$ : relate  $p$  to its children  $u, v$  via  $\text{check}(p, u, v)$



- Two chances for verification:
  - Top down using  $p, \bar{u}, v$
  - Bottom up using  $\bar{u}, v, \bar{p}$

# Binary Trees - general ideas

- **Goal:** for all internal nodes  $p$ : relate  $p$  to its children  $u, v$  via  $\text{check}(p, u, v)$



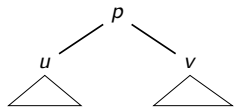
- Two chances for verification:
  - Top down using  $p, \bar{u}, v$
  - Bottom up using  $\bar{u}, v, \bar{p}$
- $\bar{u}, v$  is used for verification in any case

**Strategies:** store either  $p$  until  $\bar{u}v$  arrives, or throw  $p$  away and store  $\bar{u}v$  until  $\bar{p}$  arrives



# Stack Algorithm

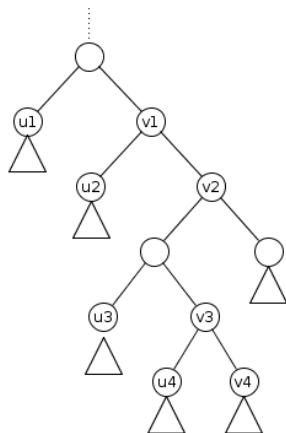
**1st idea:** Start with stack algorithm doing bottom-up verifications



subtree of  $u$     subtree of  $v$

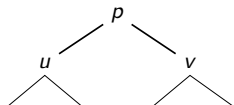
$\dots pu \dots \bar{u}v \dots \bar{v}\bar{p} \dots$

- Ignore opening tags of parent nodes



# Stack Algorithm

**1st idea:** Start with stack algorithm doing bottom-up verifications

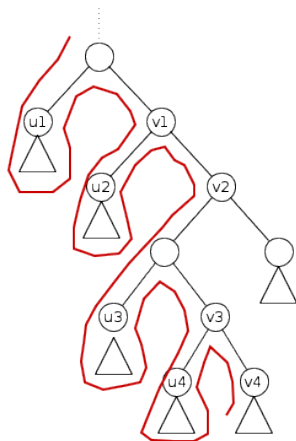


subtree of  $u$     subtree of  $v$

$\dots pu \dots \bar{u}v \dots \bar{v}p \dots$

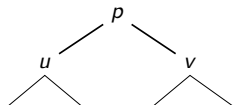
- Ignore opening tags of parent nodes
- Push children information on a stack:

$\vdots$
$\bar{u}_4, v_4$
$\bar{u}_3, v_3$
$\bar{u}_2, v_2$
$\bar{u}_1, v_1$
$\vdots$



# Stack Algorithm

**1st idea:** Start with stack algorithm doing bottom-up verifications



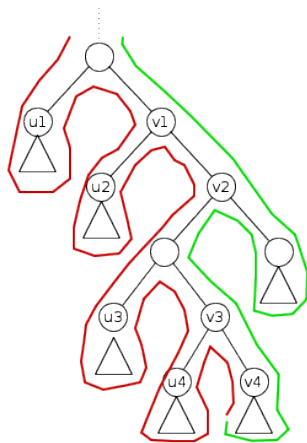
subtree of  $u$     subtree of  $v$

$\dots pu \dots \bar{u}v \dots \bar{v}p \dots$

- Ignore opening tags of parent nodes
- Push children information on a stack:

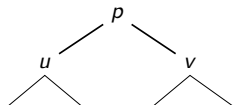
$\vdots$
$\bar{u}_4, v_4$
$\bar{u}_3, v_3$
$\bar{u}_2, v_2$
$\bar{u}_1, v_1$
$\vdots$

- Verify when going up



# Stack Algorithm

**1st idea:** Start with stack algorithm doing bottom-up verifications



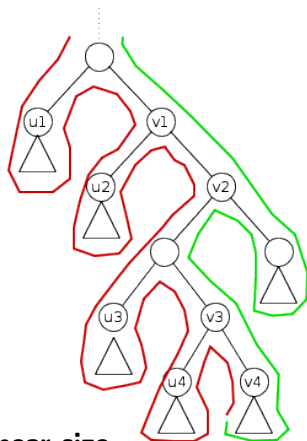
subtree of  $u$     subtree of  $v$

$\dots pu \dots \bar{u}v \dots \bar{v}p \dots$

- Ignore opening tags of parent nodes
- Push children information on a stack:

$\vdots$
$\bar{u}_4, v_4$
$\bar{u}_3, v_3$
$\bar{u}_2, v_2$
$\bar{u}_1, v_1$
$\vdots$

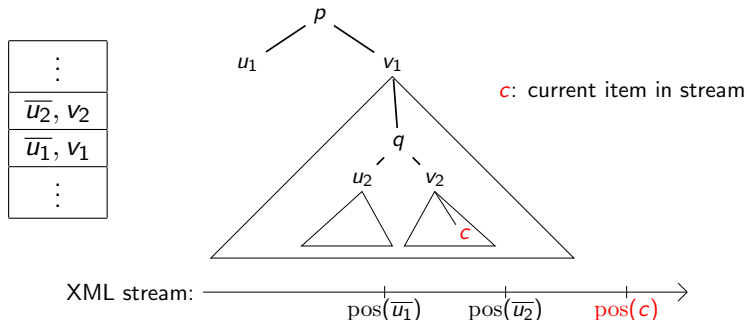
- Verify when going up



- **Stack of linear size**
- **Verification of all nodes**

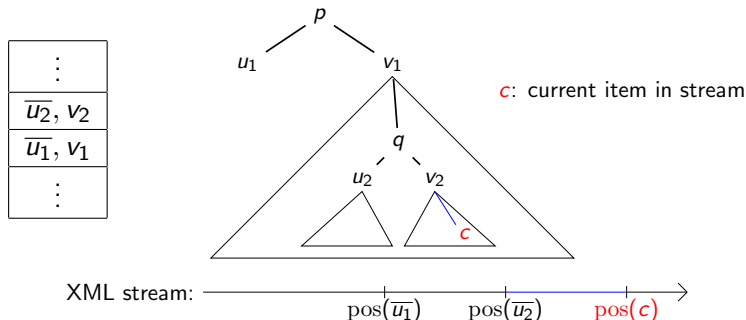
# Verifying nodes with larger left subtree

**2nd idea:** Reduce stack to  $\log(n)$  elements: remove children  $\bar{u}v$  from stack whose parents' node has a smaller left subtree than its right subtree



# Verifying nodes with larger left subtree

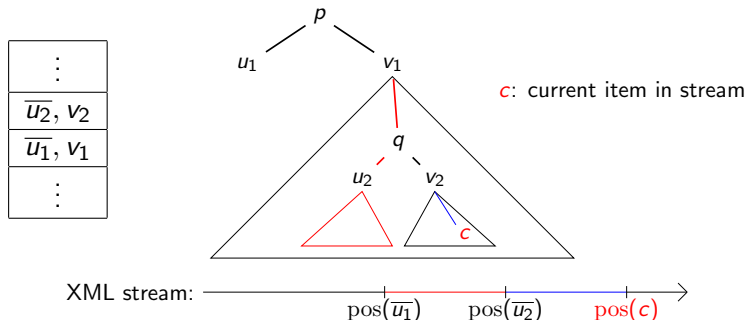
**2nd idea:** Reduce stack to  $\log(n)$  elements: remove children  $\bar{u}v$  from stack whose parents' node has a smaller left subtree than its right subtree



- $\text{pos}(c) - \text{pos}(\bar{u}_2) \leq \text{size of right subtree of } q$

# Verifying nodes with larger left subtree

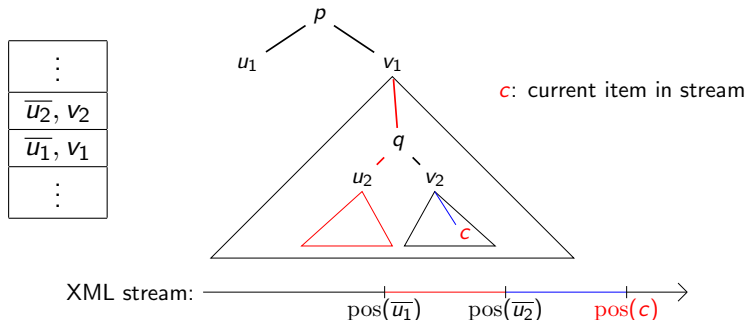
**2nd idea:** Reduce stack to  $\log(n)$  elements: remove children  $\bar{u}v$  from stack whose parents' node has a smaller left subtree than its right subtree



- $\text{pos}(c) - \text{pos}(\bar{u}_2) \leq \text{size of right subtree of } q$
- $\text{pos}(\bar{u}_2) - \text{pos}(\bar{u}_1) \geq \text{size of left subtree of } q$

# Verifying nodes with larger left subtree

**2nd idea:** Reduce stack to  $\log(n)$  elements: remove children  $\bar{u}v$  from stack whose parents' node has a smaller left subtree than its right subtree



- $\text{pos}(c) - \text{pos}(\bar{u}_2) \leq \text{size of right subtree of } q$
- $\text{pos}(\bar{u}_2) - \text{pos}(\bar{u}_1) \geq \text{size of left subtree of } q$

**Deletion rule:** delete if  $\text{pos}(c) - \text{pos}(\bar{u}_2) > \text{pos}(\bar{u}_2) - \text{pos}(\bar{u}_1)$





## Verifying nodes with larger left subtree (3)

**Lemma:** Stack is of size at most  $\log(N)$ .

**Proof:**

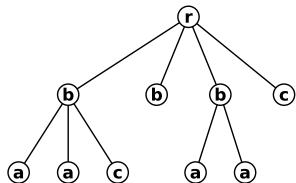
- **Deletion rule:**  $\text{pos}(c) - \text{pos}(\bar{u}_i) > \text{pos}(\bar{u}_i) - \text{pos}(\bar{u}_{i-1})$
- $u_i v_i$  **remains on stack:**  $\Rightarrow$  deletion rule does not apply

$\bar{u}_{max}, v_{max}$
$\vdots$
$\bar{u}_2, v_2$
$\bar{u}_1, v_1$

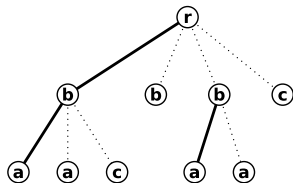
$$\begin{aligned}\text{pos}(\bar{u}_i) &> \frac{\text{pos}(c) + \text{pos}(\bar{u}_{i-1})}{2} \\ \text{pos}(\bar{u}_2) &\geq \lceil \frac{\text{pos}(c)}{2} \rceil, \\ \text{pos}(\bar{u}_3) &\geq \lceil \frac{3\text{pos}(c)}{4} \rceil, \\ &\dots\end{aligned}$$

- This leaves only space for  $\log(\text{pos}(c))$  elements □

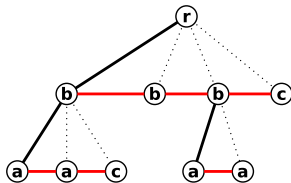
# First-Child-Next-Sibling encoding



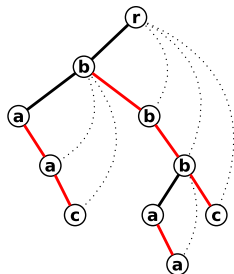
1. Original document



2. Keep edges to first children

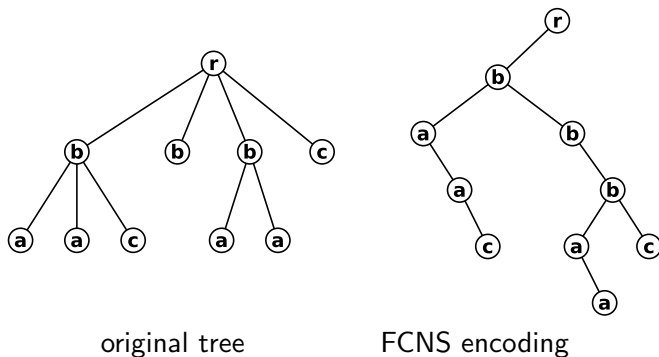


3. Insert edges connecting children



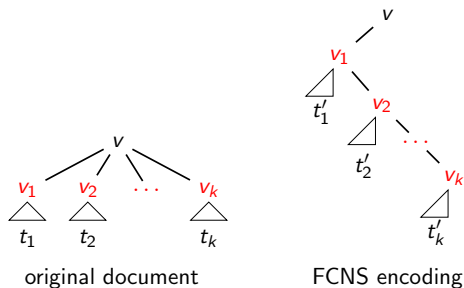
4. FCNS encoding

# First-Child-Next-Sibling encoding

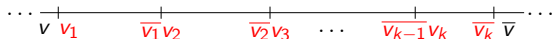


- **Transformation:** For each node in original document:
  - **First child:** becomes *left child* of that node
  - **Next Sibling:** becomes *right child* of that node
- **Annotation:** tags in FCNS encoding are annotated *left/right*

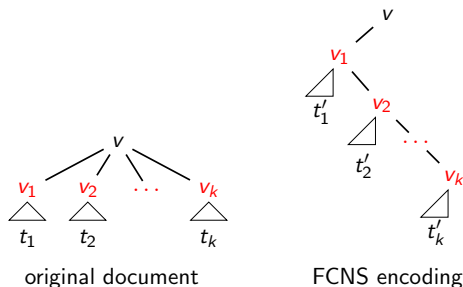
# Validation is easier given the FCNS encoding



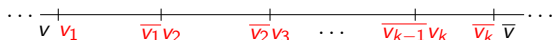
- **Original document:** tags of children of  $v$  scattered



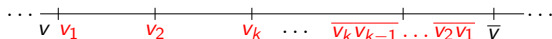
# Validation is easier given the FCNS encoding



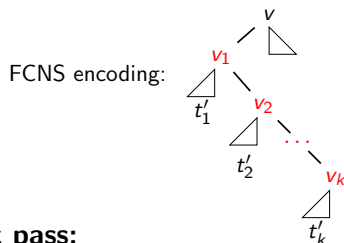
- **Original document:** tags of children of  $v$  scattered



- **FCNS encoding:**  $\overline{v_k v_{k-1}} \dots \overline{v_2 v_1}$  appears as substring



# Reusing binary tree validation algorithm



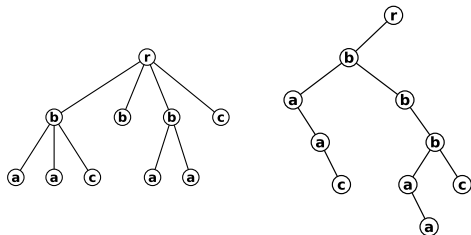
- **Left-to-Right pass:**

- compress subsequence  $\overline{v_k} \dots \overline{v_1}$  via an automaton  $A_L$  constructed from initial DTD into a state
- annotate  $\overline{v_1}$  with that state
- binary tree validation algorithm relates state to label of parent

- **Right-to-Left pass:**

- compress subsequence  $v_1 \dots v_k$  via an automaton  $A_R$  constructed from initial DTD into a state
- if binary tree algorithm pushed  $v_1$  onto stack, annotate stack element by this state
- binary tree validation algorithm relates state to label of parent

# Computing the FCNS encoding

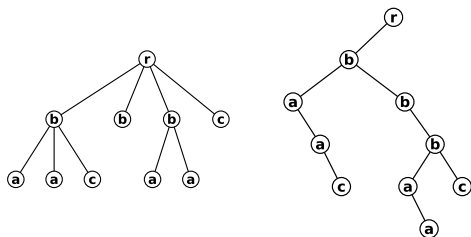


*rbāāāāc̄c̄bb̄bb̄baāāāb̄c̄c̄r rbaacc̄aabb̄baāāāc̄c̄bb̄b̄r*

- **Computing the FCNS encoding:** reordering of XML tags and annotation
- **Algorithm:**
  - 1 **compute sequence of opening tags with annotations**  
sequences of opening tags coincide
  - 2 **compute sequence of closing tags with annotations**  
start with sequence of opening tags, interpret them as closing tags, and reorder them via a modified merge sort
  - 3 **merge these sequences**



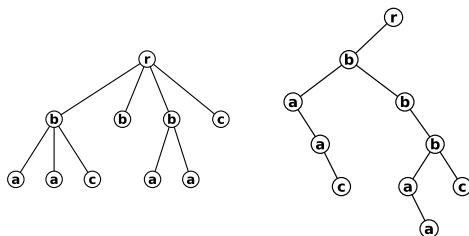
# Computing the FCNS encoding



*rbāāāc̄bb̄bb̄āāāb̄c̄c̄r rbaac̄c̄aabb̄āāāc̄c̄bb̄b̄r*

- **Computing the FCNS encoding:** reordering of XML tags and annotation
- **Algorithm:**
  - 1 **compute sequence of opening tags with annotations**  
sequences of opening tags coincide
  - 2 **compute sequence of closing tags with annotations**  
start with sequence of opening tags, interpret them as closing tags, and reorder them via a modified merge sort
  - 3 **merge these sequences**

# Computing the FCNS encoding



$rba\bar{a}\bar{a}\bar{c}\bar{b}\bar{b}\bar{b}ba\bar{a}\bar{a}\bar{b}\bar{c}\bar{c}\bar{r}$   $rbaac\bar{c}\bar{a}\bar{a}\bar{b}bba\bar{a}\bar{a}\bar{c}\bar{c}\bar{b}\bar{b}\bar{b}\bar{r}$

- **Computing the FCNS encoding:** reordering of XML tags and annotation
- **Algorithm:**
  - 1 **compute sequence of opening tags with annotations**  
sequences of opening tags coincide
  - 2 **compute sequence of closing tags with annotations**  
start with sequence of opening tags, interpret them as closing tags, and reorder them via a modified merge sort
  - 3 **merge these sequences**

## We have:

- One pass,  $O(\sqrt{N \log N})$  space for two-ranked trees
- Two bidirectional passes,  $O(\log^2 N)$  space for two-ranked trees
- $O(\log N)$  passes, 3 aux. streams,  $O(\log^2 N)$  space for arbitrary trees

## Open Problems:

- Lower bound: optimality of one pass algorithm for binary trees
- Lower bound:  $\Omega(\log(N))$  passes are required for unranked trees when using sublinear space and a constant number of auxiliary streams
- Other membership problems:  $\text{DYCK}(k) \subset \text{Visibly Pushdown languages} \subset \text{deterministic context free languages} \subset \text{context free languages}$