

Querying and storing XML

Week 7
Typechecking and Static Analysis
March 5-8, 2013

1

Why type-check XML transformations?

- Goal: Check that valid input -> valid output
 - e.g. check that XSLT/XQuery produces valid HTML
- Can we prove this once and for all?
 - rather than revalidating after every run?
- Typechecking can also support other goals:
 - optimization/static analysis
 - identifying silly bugs ("dead" code that never matches anything; missing cases)

Qsx

March 5-8, 2013

3

Regular expression types

Qsx

March 5-8, 2013

2

Example: XSLT

- Turns a sequence of People into a table
- HTML Tables have to have at least one row!

```
<xsl:output method="html"/>
<xsl:template match="/person">
  <html>
    <head><title>Phone Book</title></head>
    <body>
      <table frame="box" rules="all">
        <xsl:apply-templates/>
      </table>
    </body>
  </html>
</xsl:template>

<xsl:template match="record">
  <tr>
    <td><xsl:copy-of select="name"/></td>
    <td><xsl:copy-of select="phone"/></td>
  </tr>
</xsl:template>
```

Qsx

March 5-8, 2013

4

Influential approach: XDuce

- A functional language for transforming XML based on pattern matching

```
type Person = person[name[String],  
                    email[String]*,  
                    phone[String]?]  
  
fun person2row(var p as Person) : Tr =  
  match p with  
  person[name[var name],  
         email[var email]] ->  
  tr[td[name],td[email]]  
  | ...
```

Qsx

March 5-8, 2013

5

Regular expression types

- $T ::= \text{String} \mid X \mid a[T] \mid T, T \mid T|T \mid T^*$
 - $a[T]$ means "element labeled a with content T "
 - X is a *type variable*
- Can model DTDs, XML Schemas as special case
 - at least as far as element structure goes
 - does not model attributes, or XML Schema interleaving

Qsx

March 5-8, 2013

7

XDuce, continued

```
fun people2rows(var xs as Person+) : Tr+ =  
  match xs with  
    var p as Person, var ps as Person*  
    -> person2row p, people2rows ps  
  | var p as Person -> person2row p
```

```
fun people2table(var xs as Person+) : Table =  
  table[people2rows xs]
```

Qsx

March 5-8, 2013

6

Recursive type definitions

- Consider type definitions

```
type X1 = T1  
type X2 = T2 ...
```
- Definitions may be recursive, but recursion must be "guarded" by element tag
- e.g. type $X = X, X$ not allowed

Qsx

March 5-8, 2013

8

Meaning of types

- Suppose $E(X)$ is the definition of X for each type variable X (i.e. `type X = E(X)` declared)
- Define $\llbracket T \rrbracket_E$ as follows:
 - $\llbracket X \rrbracket_E = \llbracket E(X) \rrbracket_E$
 - $\llbracket \text{String} \rrbracket_E = \Sigma^*$
 - $\llbracket a[T] \rrbracket_E = \{a[v] \mid v \in \llbracket T \rrbracket_E\}$
 - $\llbracket T_1, T_2 \rrbracket_E = \{v_1, v_2 \mid v_1 \in \llbracket T_1 \rrbracket_E, v_2 \in \llbracket T_2 \rrbracket_E\}$
 - $\llbracket T_1 | T_2 \rrbracket_E = \llbracket T_1 \rrbracket_E \cup \llbracket T_2 \rrbracket_E$
 - $\llbracket T^* \rrbracket_E = \{v_1, \dots, v_n \mid v_1, \dots, v_n \in \llbracket T \rrbracket_E\}$

QSX

March 5-8, 2013

9

DTDs as a special case

- DTD rules:
 - `a → b, c*, (d|e)?, f`
 - `b → c, d ...`
- becomes:
 - `type A = a[B, C*, (D|E)?, F]`
 - `type B = b[C, D] ...`
- (just introduce new type `elt` for each element name `elt`)
 - Same idea works for XML Schemas

QSX

March 5-8, 2013

10

Complications

- Typechecking undecidable for any Turing-complete language
 - and for many simpler ones
- Schema languages for XML use regular languages/tree automata
 - decision problems typically EXPTIME-complete
- Nevertheless, efficient-in-practice techniques can be developed

QSX

March 5-8, 2013

11

Other checks

- Typing: All results are contained in declared types
- Exhaustiveness: every possible input matches *some* pattern
 - `match (x:Person) with`
 - `person[name[var n]] -> ...`
- Irredundancy: no pattern is "vacuous"
 - `match (x:Person) with`
 - `person[name[var n], var es as Email*] -> ...`
 - `| person[name[var n], email[var e]] -> ...`
- All can be reduced to **subtype** checks

QSX

March 5-8, 2013

12

Subtyping as inclusion

- We can define a semantic subtyping relation:
 - $T <: T' \Leftrightarrow \llbracket T \rrbracket_E \subseteq \llbracket T' \rrbracket_E$
- i.e. every tree matching T also matches T' .
- How can we check this?
- One solution: *tree automata*

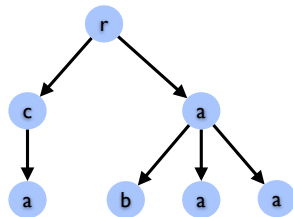
QSX

13

March 5-8, 2013

XML trees as binary trees

- First-child, next-sibling encoding



- Leaves are "nil" / ()

QSX

15

March 5-8, 2013

Type inclusion & tree automata

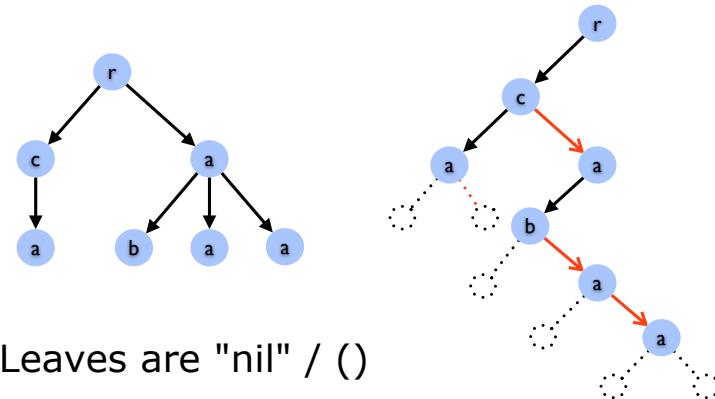
QSX

14

March 5-8, 2013

XML trees as binary trees

- First-child, next-sibling encoding



- Leaves are "nil" / ()

QSX

15

March 5-8, 2013

Binary tree automata (bottom-up)

- $M = (Q, \Sigma, \delta, q_0, F)$ is a (bottom-up) binary tree automaton if:
 - $q_0 \in Q$ -- initial state (leaf nodes)
 - $F \subseteq Q$ -- set of final states
 - $\delta : \Sigma \times Q \times Q \rightarrow Q$ -- transition function

Binary tree automata (top-down)

- $M = (Q, \Sigma, \delta, q_0, F)$ is a (top-down) binary tree automaton if:
 - $q_0 \in Q$ -- initial state (root node)
 - $F \subseteq Q$ -- set of final states (for all leaf nodes)
 - $\delta : Q \times \Sigma \rightarrow Q \times Q$ -- transition function
- Deterministic top-down TA are strictly less expressive than bottom-up TA
 - however, nondeterministic versions equivalent

DTDs and ambiguity

- Recall DTDs have to be "unambiguous"
- Example: NOT $(a+b)+(a+c)$
 - but $a+b+c$ is equivalent
- NOT $(a+b)^*a$
 - but $(b^*a)^+$ is equivalent
- These restrictions (and similar ones for XML Schemas) lead to deterministic, top-down (tree) automata.

Translating types

Translating types

- First give each subexpression its own type name

```
type Person = person[name[String],  
                    email[String]*,  
                    phone[String]?]
```

Translating types

- First give each subexpression its own type name

```
type Person = person[name[String],  
                    email[String]*,  
                    phone[String]?]
```

```
type Person = person[Name,  
                    Emails,  
                    PhoneOpt]  
type Name      = name[String]  
type Emails    = email[String], Emails | ()  
type PhoneOpt  = phone[String] | ()
```

To binary form

- Next reorganize into binary steps

```
type Person = person[Q1], Empty  
type Q1 = name[String], Q2  
         | name[String], Q3  
type Q2 = email[String], Q2  
         | email[String], Q3  
type Q3 = phone[String], Empty | ()  
type Empty = ()
```

To tree automaton

```
type Person = person[Q1], Empty  
type Q1 = name[String], Q2  
         | name[String], Q3  
type Q2 = email[String], Q2  
         | email[String], Q3  
type Q3 = phone[String], Empty | ()  
type Empty = ()
```

```
Q = {Person, Q1, Q2, Q3, String, Empty}  
Σ = {person, name, email, phone, string}  
q0 = Person  
F = {Empty, Q3, String}
```

To tree automaton

```
type Person = person[Q1],Empty
type Q1 = name[String], Q2
      | name[String], Q3
type Q2 = email[String],Q2
      | email[String],Q3
type Q3 = phone[String],Empty | ()
type Empty = ()
```

```
Q = {Person,Q1,Q2,Q3,String}
Σ = {person,name,email,phone}
q0 = Person
F = {Empty,Q3,String}
```

q ₁	q ₂	a	δ(q ₁ ,q ₂ ,a)
Q1	Empty	person	Person
String	Q2	name	Q1
String	Q3	name	Q1
String	Q2	email	Q2
String	Q3	email	Q3
String	Empty	phone	Q3

Qsx

March 5-8, 2013

21

Checking containment

- Like sequential case, tree automata / languages closed under union, intersection, negation
 - some of the constructions incur exponential blowup
- Containment testing is EXPTIME-complete
 - Naive algorithm: test $L(M) \subseteq L(N)$ by constructing automaton M' for $L(M) \cap L(\bar{N})$, testing emptiness
 - Does not perform well on problems encountered in programming
 - Hosoya, Vouillon & Pierce [2005] give practical algorithm aimed at XDuce typing problems

Qsx

March 5-8, 2013

22

Types for XQuery

Qsx

March 5-8, 2013

23

Motivation

[Collazzo et al. 2006]

- For XDuce, want to ensure valid result
 - find silly bugs in pattern matching
 - ensure exhaustiveness & irredundancy
- For XQuery, still want result validity, but different issues
 - functions, but no pattern matching (as such)
 - but does have for/iteration & path expressions
 - parts of queries can have "silly" bugs

Qsx

March 5-8, 2013

24

Example

- Schema (we'll use XQuery notation):

```
type Person = person[name[String],
                      email[String]*,
                      phone[String]?]
```

- Query: return all name/phone of people with email

```
for $y in $x/person[emial]
return $y/name,$y/phone
```

- **Silly bug:** "emial" misspelled; path will never select anything (in this type anyway)

Q SX

25

March 5-8, 2013

Typing rules for XQuery: Basics

$$\frac{\bar{x}:\alpha \in \Gamma}{\Gamma \vdash \bar{x}:\alpha} \quad \frac{x:\tau \in \Gamma}{\Gamma \vdash x:\tau} \quad \frac{w \in \Sigma^*}{\Gamma \vdash w:\mathbf{string}} \quad \frac{b \in \mathbf{Bool}}{\Gamma \vdash b:\mathbf{bool}}$$

$$\frac{}{\Gamma \vdash () : ()} \quad \frac{\Gamma \vdash e:\tau}{\Gamma \vdash n[e]:n[\tau]} \quad \frac{\Gamma \vdash e:\tau \quad \Gamma \vdash e':\tau'}{\Gamma \vdash e, e':\tau, \tau'}$$

Q SX

27

March 5-8, 2013

XQuery core language: μ XQ

$$e ::= () \mid e, e' \mid n[e] \mid w \mid x \mid \text{let } x = e \text{ in } e' \\ \mid b \mid \text{if } c \text{ then } e \text{ else } e' \mid \bar{x} \mid \bar{x}/\text{child} \mid e :: n \mid \text{for } \bar{x} \in e \text{ return } e'$$

$$\bar{v} ::= b \mid w \mid n[v] \quad v ::= \bar{v}, v \mid ()$$

$$\Gamma ::= \cdot \mid \Gamma, x:\tau \mid \Gamma, \bar{x}:\alpha$$

Tree variables \bar{x} have single tree values \bar{v}
Forest variables x have arbitrary values v

Q SX

26

March 5-8, 2013

Typing rules for XQuery: Let, If

$$\frac{\Gamma \vdash e_1:\tau_1 \quad \Gamma, x:\tau_1 \vdash e_2:\tau_2}{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2:\tau_2}$$

$$\frac{\Gamma \vdash c:\mathbf{bool} \quad \Gamma \vdash e_1:\tau_1 \quad \Gamma \vdash e_2:\tau_2}{\Gamma \vdash \text{if } c \text{ then } e_1 \text{ else } e_2:\tau_1|\tau_2}$$

Q SX

28

March 5-8, 2013

Path steps

$$\frac{\bar{x}:n[\tau] \in \Gamma}{\Gamma \vdash \bar{x}/\text{child} : \tau} \quad \frac{\Gamma \vdash e : \tau \quad \tau :: n \Rightarrow \tau'}{\Gamma \vdash e :: n : \tau'}$$

$$\frac{\overline{n[\tau] :: n \Rightarrow n[\tau]}}{E(X) :: n \Rightarrow \tau \quad \alpha \neq n[\tau]} \quad \frac{\alpha :: n \Rightarrow ()}{X :: n \Rightarrow \tau} \quad \frac{\tau_1 :: n \Rightarrow \tau_2}{() :: n \Rightarrow ()} \quad \frac{\tau_1^* :: n \Rightarrow \tau_2^*}{\tau_1 :: n \Rightarrow \tau_1' \quad \tau_2 :: n \Rightarrow \tau_2'}$$

$$\frac{\tau_1, \tau_2 :: n \Rightarrow \tau_1', \tau_2'}{\tau_1 :: n \Rightarrow \tau_1' \quad \tau_2 :: n \Rightarrow \tau_2'}$$

$$\frac{\tau_1 | \tau_2 :: n \Rightarrow \tau_1' | \tau_2'}{\tau_1 | \tau_2 :: n \Rightarrow \tau_1' | \tau_2'}$$

Qsx

March 5-8, 2013

29

Example

```
type Doc = a[ b[d[]]*, c[e[], f[]]? ]
prime( b[d[]]*, c[e[], f[]]? ) = b[d[]]|c[e[],f[]]
quantifier( b[d[]]*, c[e[], f[]]? ) = *
```

\$doc : Doc \vdash \$doc/a/* : b[d[]]*, c[e[], f[]]?

\$doc : Doc, \$x : b[d[]] | c[e[], f[]] \vdash \$x/* : d[] | e[],f[]

\$doc : Doc \vdash for \$x in \$doc/a/* return \$x/* : (d[] | e[],f[])*

Overapproximation:

e, f, d and e, f, e, f can never happen

Qsx

March 5-8, 2013

31

Comprehensions: simple idea (what XQuery does)

statEnv I- Expr₁ : Type₁
statEnv I- VarName₁ of var expands to Variable₁
statEnv + varType(Variable₁ => prime(Type₁)) I- Expr₂ : Type₂

statEnv I- for \$VarName₁ in Expr₁ return Expr₂ : Type₂ · quantifier(Type₁)

prime(FormalItem Type)	= FormalItem Type	quantifier(FormalItem Type)	= 1
prime(empty)	= none	quantifier(empty)	= ?
prime(none)	= none	quantifier(none)	= 1
prime(Type ₁ , Type ₂)	= prime(Type ₁) prime(Type ₂)	quantifier(Type ₁ , Type ₂)	= quantifier(Type ₁), quantifier(Type ₂)
prime(Type ₁ & Type ₂)	= prime(Type ₁) prime(Type ₂)	quantifier(Type ₁ & Type ₂)	= quantifier(Type ₁), quantifier(Type ₂)
prime(Type ₁ Type ₂)	= prime(Type ₁) prime(Type ₂)	quantifier(Type ₁ Type ₂)	= quantifier(Type ₁) quantifier(Type ₂)
prime(Type?)	= prime(Type)	quantifier(Type?)	= quantifier(Type) · ?
prime(Type*)	= prime(Type)	quantifier(Type*)	= quantifier(Type) · *
prime(Type+)	= prime(Type)	quantifier(Type+)	= quantifier(Type) · +

Qsx

March 5-8, 2013

30

Comprehensions: more precisely

$$\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash \bar{x} \text{ in } \tau_1 \rightarrow e_2 : \tau_2}{\Gamma \vdash \text{for } \bar{x} \in e_1 \text{ return } e_2 : \tau_2}$$

$$\frac{\Gamma \vdash \bar{x} \text{ in } E(X) \rightarrow e : \tau}{\Gamma \vdash \bar{x} \text{ in } () \rightarrow e : ()} \quad \frac{\Gamma \vdash \bar{x} \text{ in } X \rightarrow e : \tau}{\Gamma \vdash \bar{x} \text{ in } \alpha \rightarrow e : \tau} \quad \frac{\Gamma \vdash \bar{x} \text{ in } \tau_1 \rightarrow e : \tau_2}{\Gamma \vdash \bar{x} \text{ in } \tau_1^* \rightarrow e : \tau_2^*}$$

$$\frac{\Gamma \vdash \bar{x} \text{ in } \tau_1 \rightarrow e : \tau_1' \quad \Gamma \vdash \bar{x} \text{ in } \tau_2 \rightarrow e : \tau_2'}{\Gamma \vdash \bar{x} \text{ in } \tau_1, \tau_2 \rightarrow e : \tau_1', \tau_2'}$$

$$\frac{\Gamma \vdash \bar{x} \text{ in } \tau_1 \rightarrow e : \tau_1' \quad \Gamma \vdash \bar{x} \text{ in } \tau_2 \rightarrow e : \tau_2'}{\Gamma \vdash \bar{x} \text{ in } \tau_1 | \tau_2 \rightarrow e : \tau_1' | \tau_2'}$$

Qsx

March 5-8, 2013

32

Example

```
type Doc = a[ b[d[]]*, c[e[], f[]]? ]
prime( b[d[]]*, c[e[], f[]]? ) = b[d[]]|c[e[],f[]]
quantifier( b[d[]]*, c[e[], f[]]? ) = *
```

$\$doc : Doc \vdash \$doc/a/* : b[d[]]*, c[e[], f[]]?$

$\$doc : Doc \vdash \$x \text{ in } b[d[]]* , c[e[], f[]]? \rightarrow \$x/* : d[]*, (e[],f[])?$

$\$doc : Doc \vdash \text{for } \$x \text{ in } \$doc/a/* \text{ return } \$x/* : d[]*, (e[],f[])?$

Remembers more of the regexp structure
(but overall system still approximate)

Qsx

March 5-8, 2013

33

Path errors

- Detect *path errors* by:
 - Noticing when a non-() expression has type ()
 - Tracking which parts of expression have this property
 - (propagating sets of expression locations through typing judgements)

$\frac{\text{(TYPEEMPTY)}}{WF(E; \Gamma \vdash_{\beta} () : ()); \emptyset}$	$\frac{\text{(TYPEATOMIC)}}{WF(E; \Gamma \vdash_{\beta} b : (B); \emptyset)}$
$\frac{\text{(TYPEVARLET)}}{x : T \in \Gamma \quad WF(E; \Gamma \vdash_{\beta} x : (T); \emptyset)}$	$\frac{\text{(TYPEVARFOR)}}{\bar{x} : T \in \Gamma \quad WF(E; \Gamma \vdash_{\beta} \bar{x} : (T); \emptyset)}$
$\frac{\text{(TYPEELEM)}}{E; \Gamma \vdash_{\beta} Q : (T); \mathcal{S}}$	$\frac{\text{(TYPEFOREST)}}{E; \Gamma \vdash_{\beta} Q_1 : (T_1); \mathcal{S}_1}$
$E; \Gamma \vdash_{\beta} l[Q] : (l[T]); \mathcal{S}$	$E; \Gamma \vdash_{\beta} Q_1, Q_2 : (T_1, T_2); \mathcal{S}_1 \cup \mathcal{S}_2$

Qsx

March 5-8, 2013

34

Path errors

- For expressions that visit same subexpression multiple times
 - only error if subexpression *never* produces useful output
 - take *intersection*
- e.g. for $\$x \text{ in } \$doc/person[emial] \dots$
 - not an error if $\$doc : person[email[\dots] | emial[\dots]]$

Qsx

March 5-8, 2013

35

Subtleties

- [Collazzo et al. 2006] also considers *splitting*
- E.g. considering $a[T_1|\dots|T_n]$ as $a[T_1]|\dots|a[T_n]$
- This complicates system, but can make results more precise

Qsx

March 5-8, 2013

36

Schema-based query-update independence analysis

Query-update independence

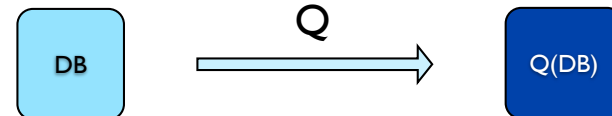


Motivation

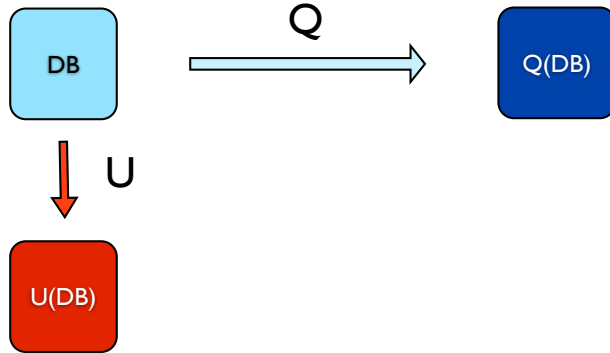
[Benedikt & Cheney, 2009]

- Suppose we want to maintain multiple (materialized) views or constraints
 - expressed by queries Q_1, Q_2, \dots
- When the database is updated
 - If we can determine (quickly) that query and update are *independent*
 - then can skip view/constraint maintenance

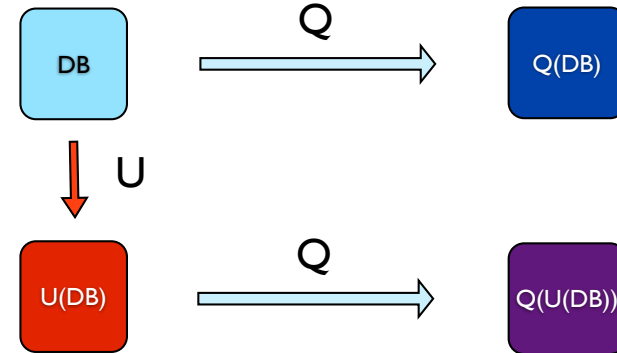
Query-update independence



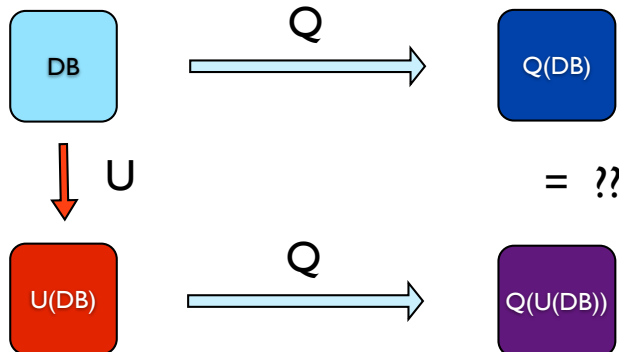
Query-update independence



Query-update independence



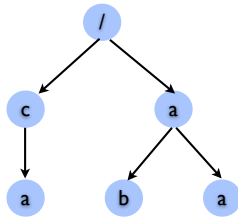
Query-update independence



Approach

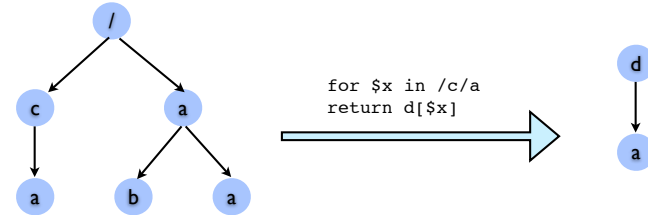
- a **static** analysis
- that safely detects **independence**
- and takes **schema** into account
- for XQuery **Update** with all XPath axes
- **fast** enough to be useful as an optimization

Independence example

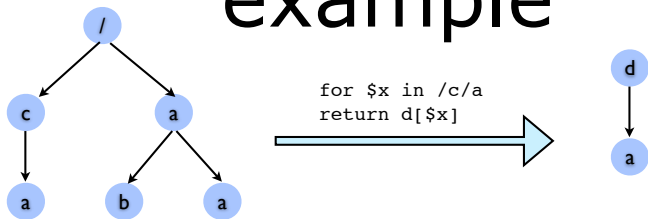


```
for $x in /c/a  
return d[$x]
```

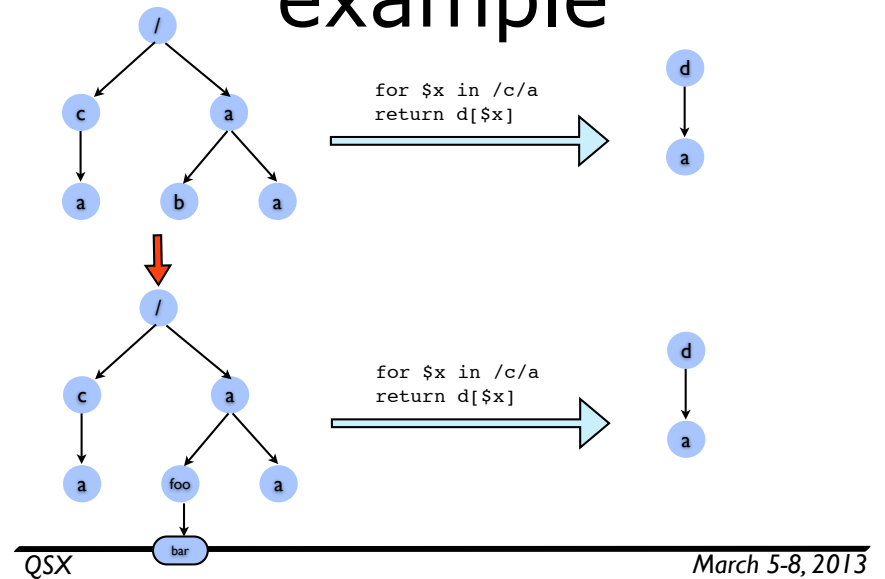
Independence example



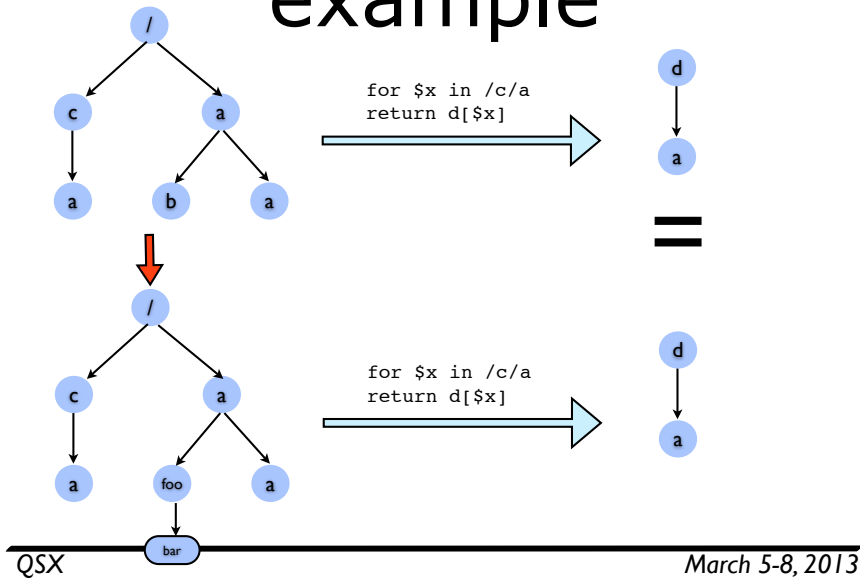
Independence example



Independence example

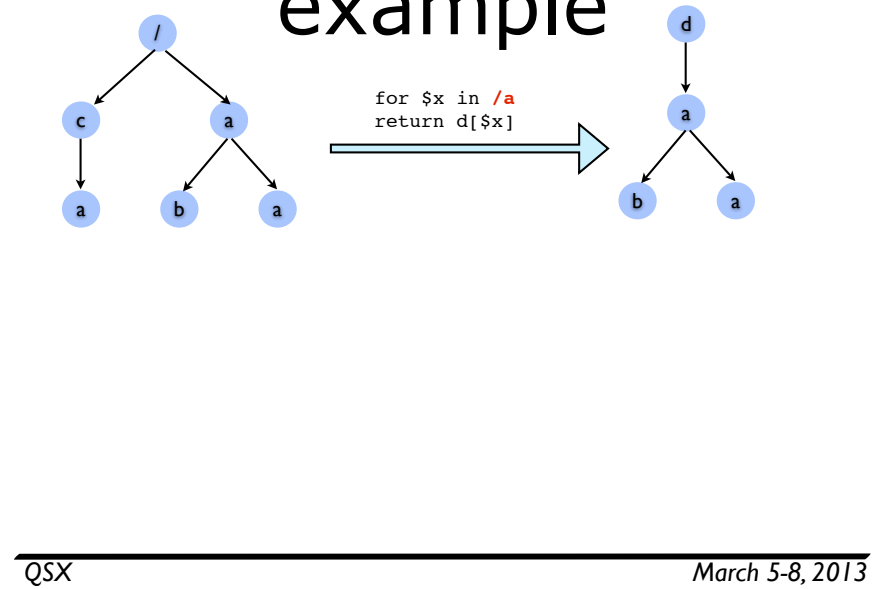


Independence example



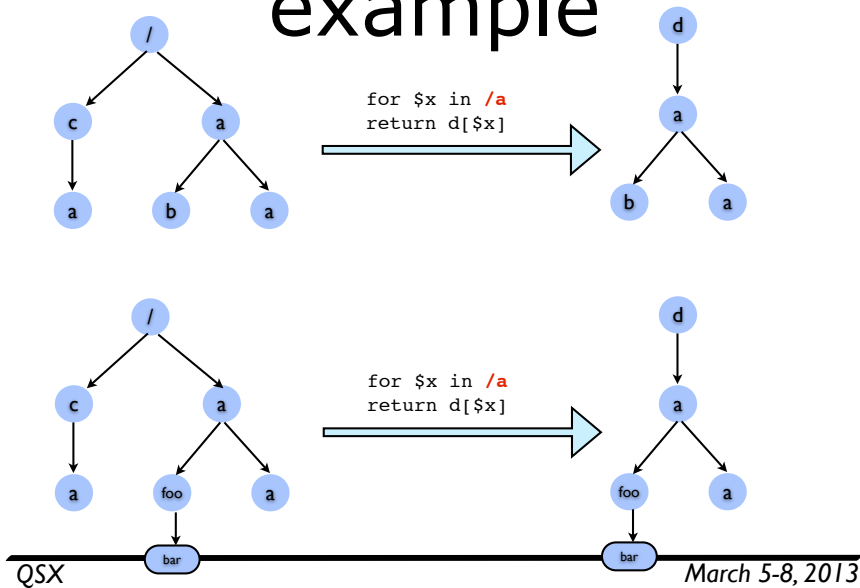
42

Independence non-example



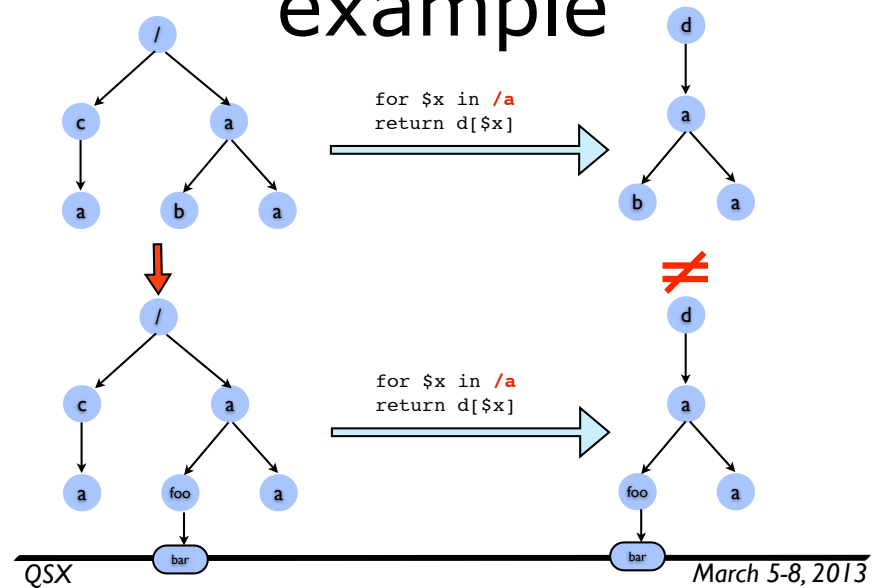
43

Independence non-example



43

Independence non-example



43

Approach

- Exploits a **schema** S that describes the input
- Statically calculate:
 - c = **Copied nodes** of Q
 - a = **Accessed Nodes** of Q
 - u = **Updated Nodes** of U
- c, a, u are sets of type names in S

Q SX

44

March 5-8, 2013

Independence test

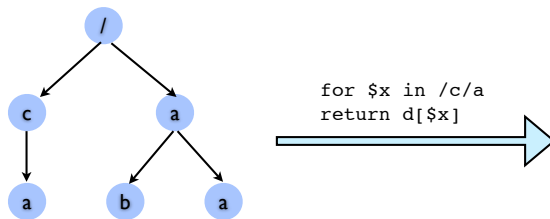
- We show that Q and I are independent modulo S if:
 - a is disjoint from u
 - that is, the update has no impact on **accessed** nodes
 - and $c//*$ is disjoint from u
 - that is, the update does not impact any copied nodes **or their descendants**

Q SX

45

March 5-8, 2013

Analysis example

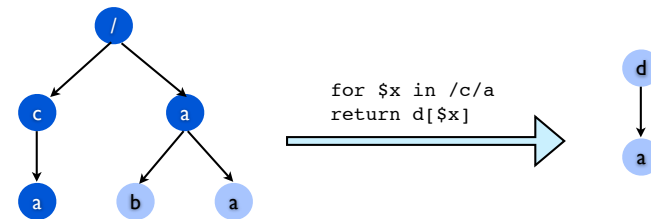


Q SX

46

March 5-8, 2013

Analysis example

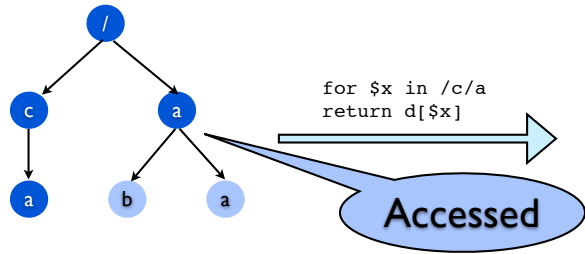


Q SX

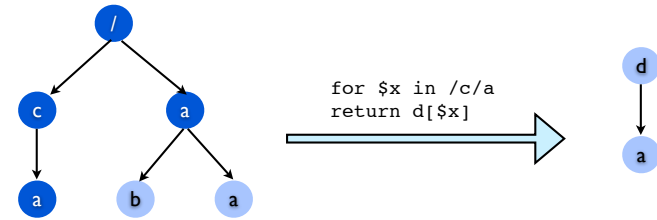
46

March 5-8, 2013

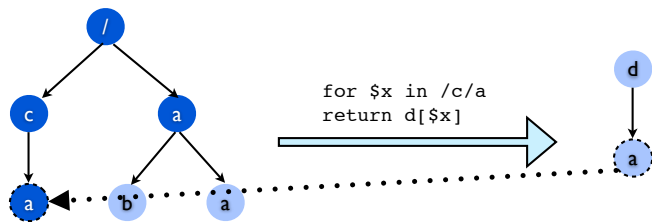
Analysis example



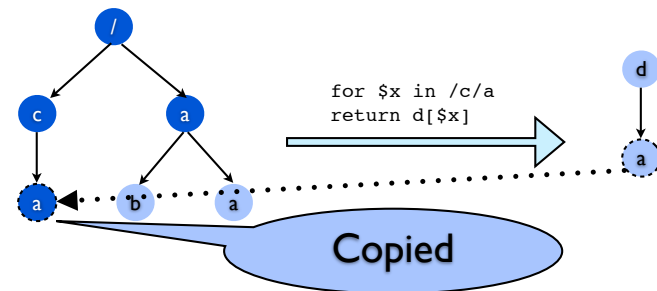
Analysis example



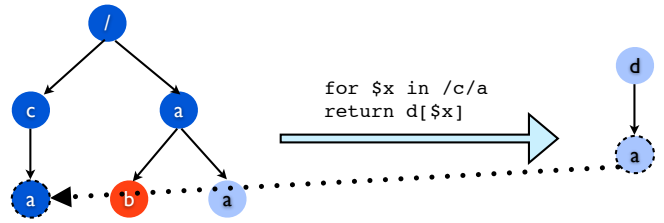
Analysis example



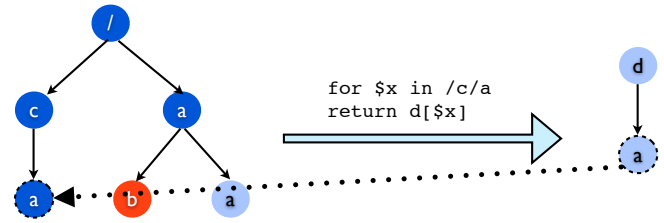
Analysis example



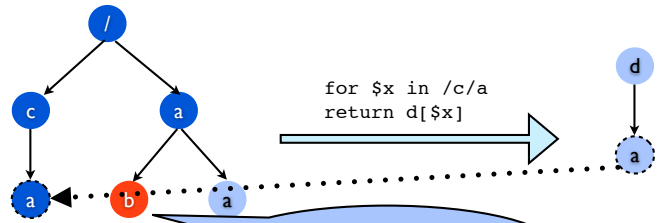
Analysis example



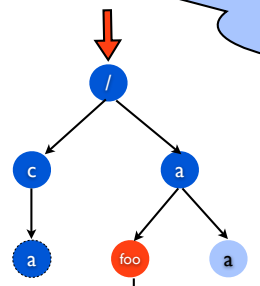
Analysis example



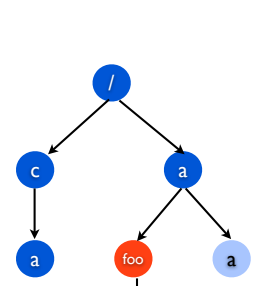
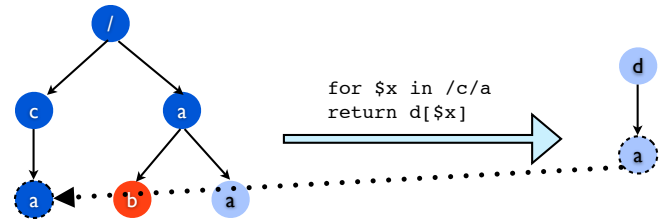
Analysis example



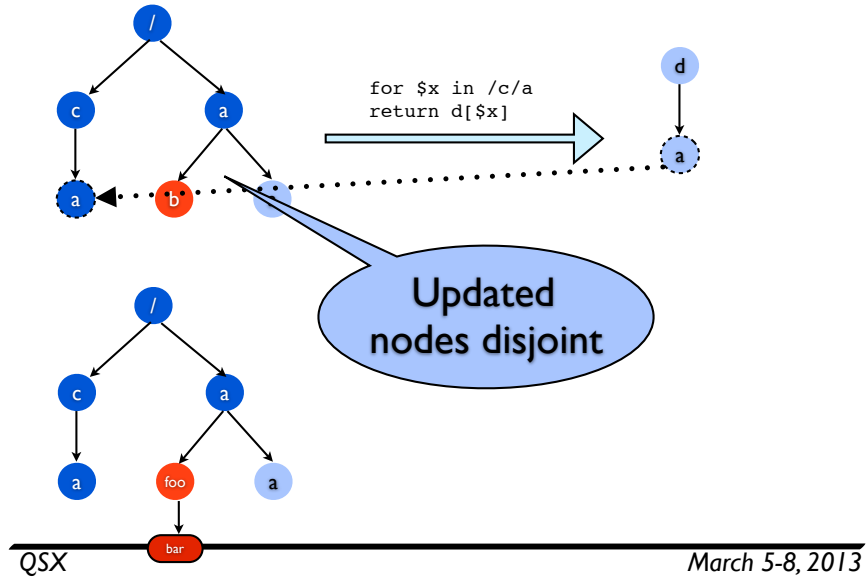
Updated



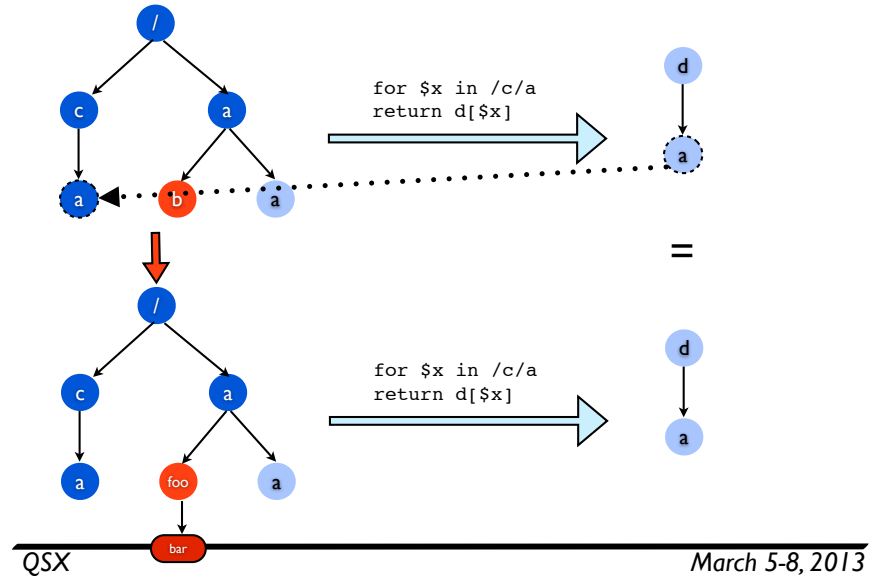
Analysis example



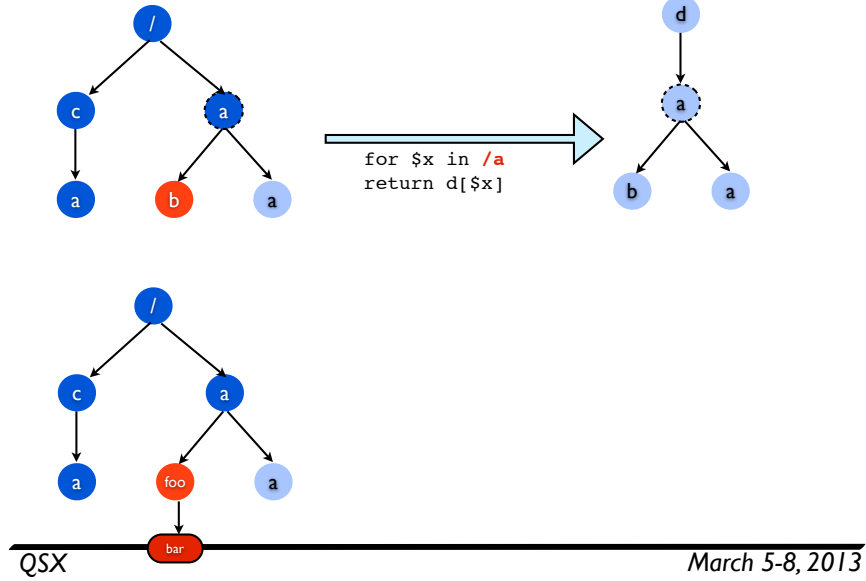
Analysis example



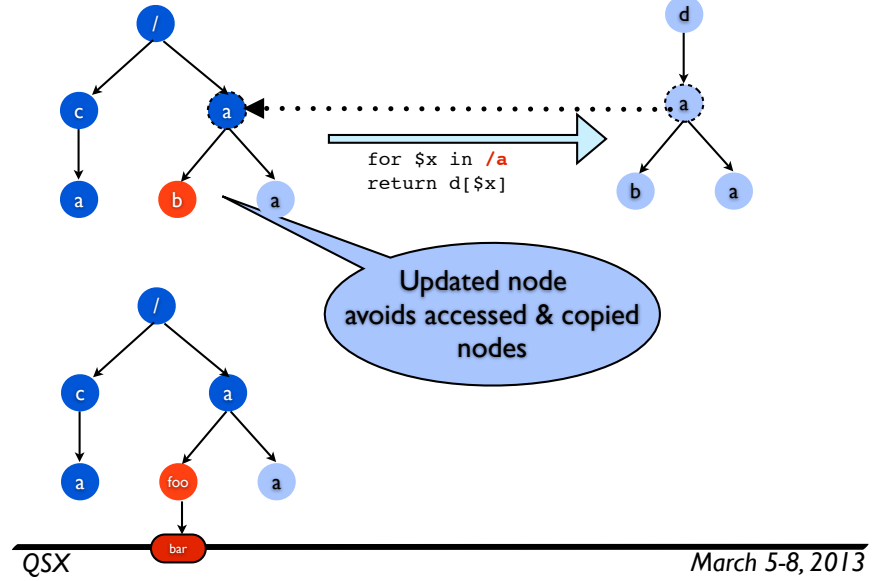
Analysis example



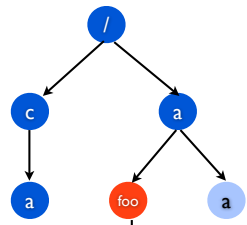
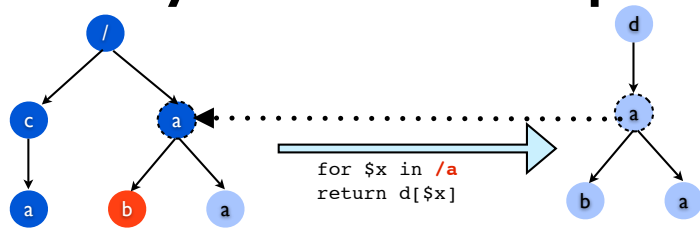
Analysis example II



Analysis example II



Analysis example II

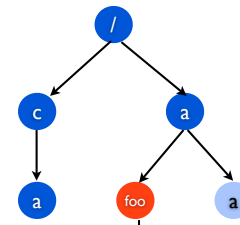
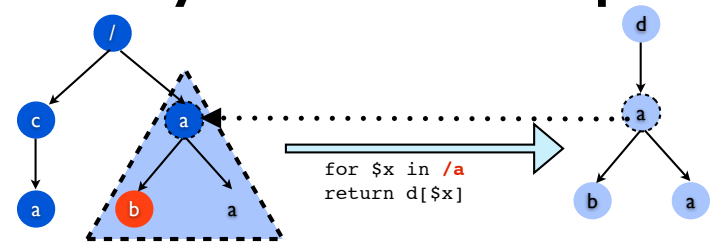


Q SX

50

March 5-8, 2013

Analysis example II

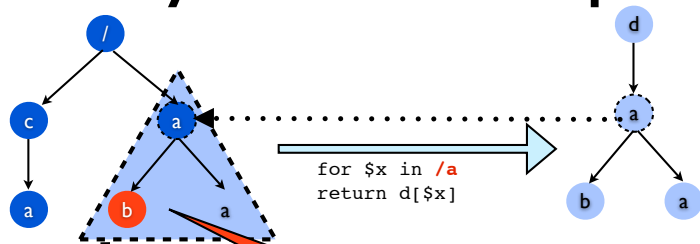


Q SX

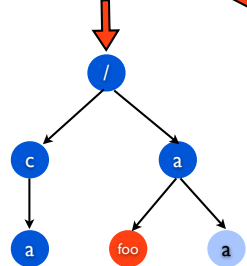
50

March 5-8, 2013

Analysis example II



Updated node is a descendant of a copied node!

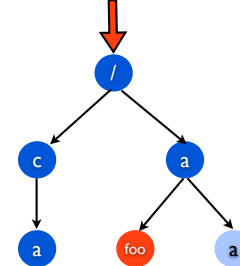
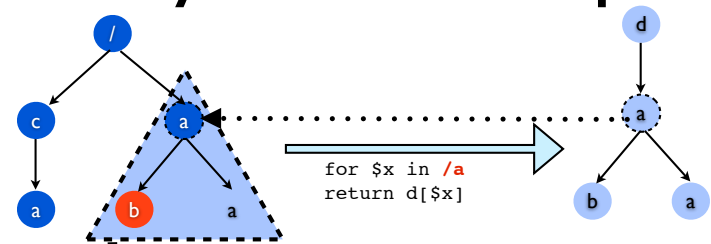


Q SX

50

March 5-8, 2013

Analysis example II

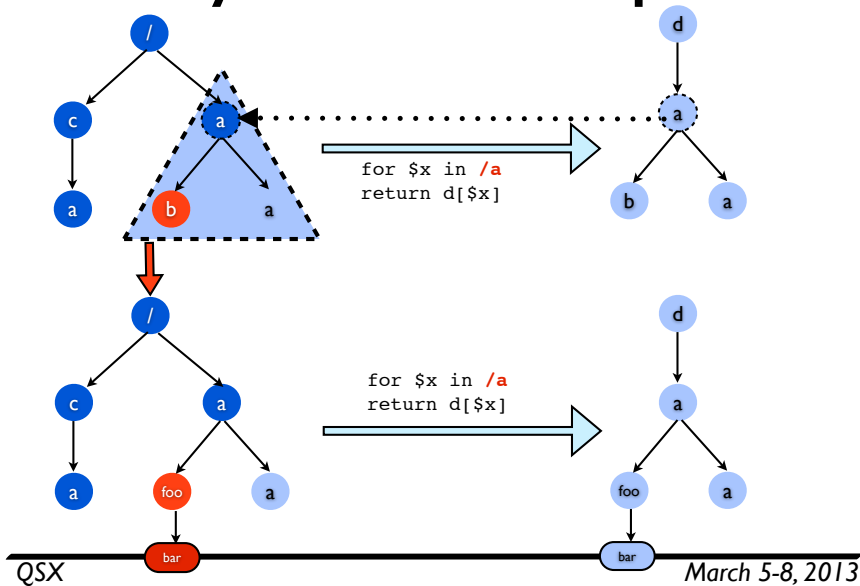


Q SX

51

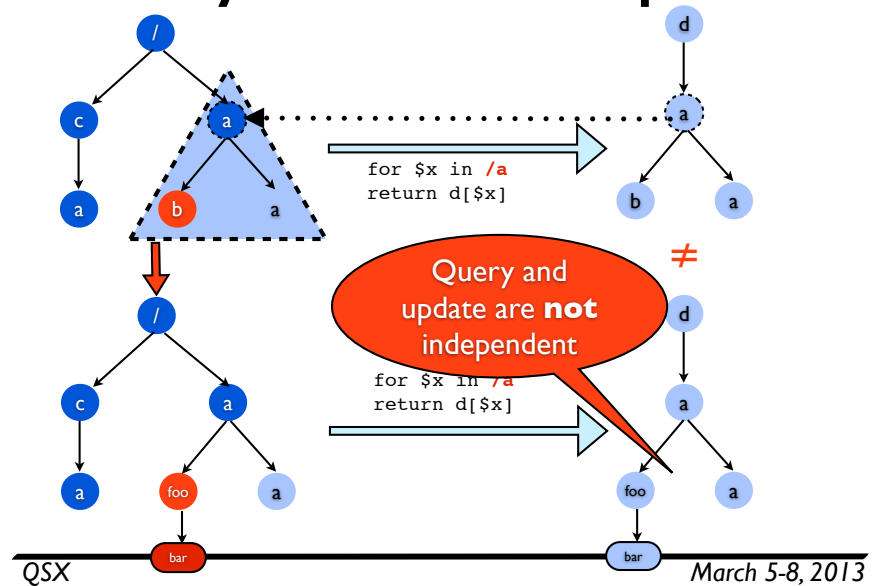
March 5-8, 2013

Analysis example II



51

Analysis example II



51

Correctness

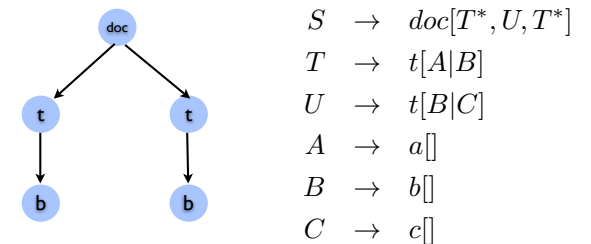
- We can use **type names** for sets of accessed, copied, updated nodes
- Symbolically evaluate query/update over **schema**
- However, there is a complication:

Aliasing

- How do we know that it is safe to use type names to refer to and change "parts" of schema?

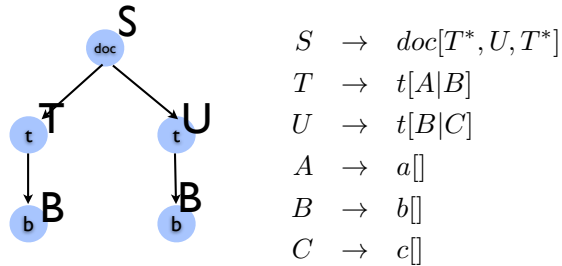
Problem

- For example:



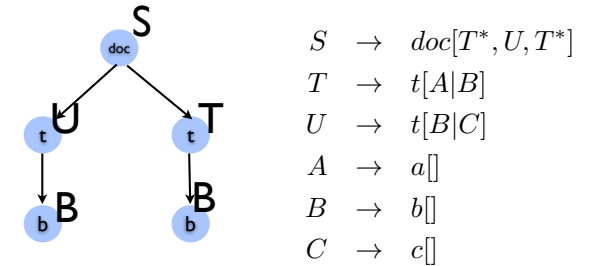
Problem

- For example:



Problem

- For example:



- has multiple typings

Solution: Closure under Aliasing

- Keep track of which type names “alias”
 - or, can refer to the same node
 - currently, based on tree language overlap test
- Close sets a , c , u under aliasing
 - example: any T could be a U , and vice versa
- Note: This may be overkill
 - unnecessary (has no effect) for DTDs
 - probably unnecessary for unambiguous XML Schema

Summary

- Independence analysis can statically determine whether view needs to be maintained
- Subsequent work explores more precise static information:
 - path-based independence analysis (“destabilizers”) - no schema required [Cheney & Benedikt VLDB 2010]
 - independence based on more precise type analysis [Bidoit-Tollu et al. VLDB 2012]
- Future work:
 - combine path- and schema-based approach?
 - what about XML views of relations?
 - can static analysis help with incremental maintenance?

Next week

- Provenance
 - Why and Where: A characterization of data provenance
 - Provenance management in curated databases
 - Annotated XML: queries and provenance