



# Backward type inference for XML queries

Hyeonseung Im<sup>a,\*</sup>, Pierre Genevès<sup>b</sup>, Nils Gesbert<sup>b</sup>, Nabil Layaïda<sup>b</sup>

<sup>a</sup> Kangwon National University, 1 Gangwondaehak-gil, Chuncheon-si, 24341 Republic of Korea

<sup>b</sup> Univ. Grenoble Alpes, CNRS, Inria, Grenoble INP, LIG, 38000 Grenoble, France



## ARTICLE INFO

### Article history:

Received 1 October 2018

Received in revised form 17 February 2020

Accepted 25 March 2020

Available online 27 March 2020

Communicated by W. Fan

### Keywords:

XPath

XQuery

Static type system

Type inference

Regular tree types

Mu-calculus

## ABSTRACT

Although XQuery is a statically typed, functional query language for XML data, some of its features such as upward and horizontal XPath axes are typed imprecisely. The main reason is that while the XQuery data model allows us to navigate upwards and between siblings from a given XML node, the type model, e.g., regular tree types, can describe only the subtree structure of the given node. To alleviate this limitation, precise forward type inference systems for XQuery were recently proposed using an extended regular type language that can describe not only a given XML node but also its context. In this paper, as a different approach, we propose a novel backward type inference system for XQuery, based on a type language extended with logical formulas. Our backward type inference system provides an exact typing result for XPath axes and a sound typing result for XQuery expressions.

© 2020 Elsevier B.V. All rights reserved.

## 1. Introduction

XQuery [1] is a statically typed, functional, World Wide Web Consortium (W3C) standard query language for XML data. Its type language is based on regular tree types (*i.e.*, regular tree languages) [2] and its static and dynamic semantics are formally defined [3]. One of the key features of XQuery is its use of XPath [4,5] to navigate and extract XML data. Although XPath navigational expressions greatly facilitate XML manipulation, they are also a main source of undesired, imprecise type inference in the XQuery formal semantics. Specifically, when upward or horizontal XPath axes such as `parent` and `following-sibling` are used, the formal semantics simply deduces the most general type (*e.g.*, `AnyElt` for `parent` and `AnyElt*` for `following-sibling` where `AnyElt` denotes the type of all XML elements), which essentially conveys no information, regardless of the type of the initial XML document. In the end, in the recent recommendations of XPath 3.0 [6] and XQuery 3.0 [7], static typing became “implementation defined” and hence optional.

The over-approximation in type inference is in particular due to the discrepancy between the XQuery data model and the type model. Specifically, in XQuery, values are *sequences of pointers* to XML tree nodes and each pointer can point anywhere in the corresponding tree. Moreover, given such a pointer, it is always possible to obtain a pointer to its parent or sibling node, thus allowing us to navigate upwards and between siblings. In clear contrast, given a pointer value, its type (*e.g.*, a regular tree type) can describe only the subtree structure to which the pointer points, but not its context, *i.e.*, part of the

\* Corresponding author.

E-mail addresses: [hsim@kangwon.ac.kr](mailto:hsim@kangwon.ac.kr) (H. Im), [pierre.geneves@cnrs.fr](mailto:pierre.geneves@cnrs.fr) (P. Genevès), [nils.gesbert@grenoble-inp.fr](mailto:nils.gesbert@grenoble-inp.fr) (N. Gesbert), [nabil.layaida@inria.fr](mailto:nabil.layaida@inria.fr) (N. Layaïda).

whole tree except the subtree pointed by the pointer value. Therefore, with this type language, only downward axes such as `child` and `desc` can be precisely typed at best (e.g., [8]).

There are two different approaches to alleviate this limitation. The first approach is to develop a typechecking algorithm based on backward type inference (also known as inverse type inference) [9–14]. Given an XQuery expression  $e$  and an expected output type  $\rho_o$ , backward type inference computes the pre-image  $\rho_i$  of  $\rho_o$  with respect to  $e$  such that it is guaranteed that for any XML document of type  $\rho_i$ ,  $e$  always produces a document of type  $\rho_o$ . Since the pre-image of a regular tree language with respect to a macro tree transducer (MTT) is also regular [15], MTTs and their variants have often been used as a model of XML transformations in the context of backward type inference [11,12,14]. Although exact typechecking can be done with backward type inference, its complexity is hyper-exponential (i.e., a stack of exponentials) [10,12,16]. To our knowledge, both practical and exact backward type inference for general XML transformations exploiting backward axes such as `parent` and `anc` has not been reported yet.

In contrast, the second approach is to develop an approximate but practical forward type inference system by using a refined type language that can describe not only XML nodes but also their contexts. For example, Castagna et al. [17] extend regular tree types with zipper data structures [18] and propose a precise type system for XQuery 3.0 (which is not exact but deduces a more precise type than the most general type such as `AnyElt` whenever possible). Their type system supports all navigational XQuery expressions including type and value case analysis and higher-order functions. Genevès and Gesbert [19] also develop a precise type system for XQuery by combining regular tree types with modal logic formulas [20]. By encoding context information using modal formulas, their type system also deduces precise types for backward axes as well as forward axes. However, none of [17] and [19] provide exact typing for XPath axes. Moreover, although practical implementation is feasible, forward type inference cannot be exact if it infers a regular tree type for admissible outputs since a general transformation does not preserve regularity. (Typechecking based on forward type inference can be exact though if it infers a more expressive type such as a context-free tree grammar [21] or a higher-order recursion scheme [22] and checks inclusion against the output type specified by a regular tree grammar.)

In this paper, we revisit a problem of backward type inference for XML queries. In particular, we develop a novel XQuery source language type system using the refined type language proposed in [19]. While tree transducers can be used as an intermediate language for XQuery, having a source language type system in itself is useful as it is usually easier to understand. Moreover, by building a backward type inference system on the XQuery syntax and the existing type language, it would be possible to combine it with forward type inference, for example, in order to develop a more precise and practical bidirectional typechecking algorithm. Thus, this work can be considered as a stepping stone towards such bidirectional type systems.

To develop a backward type inference system, we first define the syntax and semantics of an XQuery core by representing XML nodes as *focused trees* [20] (Section 2). A focused tree is a variant of zipper data structures [18], which describes a whole tree “seen” from a given internal node, that is, a subtree and its context. As focused trees support functional navigation in any direction from a given tree node, we can simplify the semantics of the XQuery core, without resorting to an external store for node pointers as in the XQuery formal semantics. With focused trees, our semantics is a straightforward extension of the one given in [8] with non-downward XPath axes.

As for our type language, we use formula-enriched sequence types [19], which combine the usual regular tree types with tree logic formulas [20] to describe both a tree node and its context (Section 3). Then, using formula-enriched sequence types, we define an exact backward type inference system for XPath axes (Section 4). That is, given an XPath axis and an output type  $\rho$ , if our inference system infers an input type  $\rho'$ , the result of evaluating the axis is of type  $\rho$  if and only if an input focused tree is of type  $\rho'$ . Then, building on the inference rules for XPath axes, we define a sound backward type inference system for the XQuery core (Section 5). In the presence of an arbitrary `FOR`-expression with a formula-enriched sequence type as an output type, both practical and exact typing is nontrivial or even may be infeasible, and therefore we introduce an approximation.

We summarize the main contributions as follows:

- We formulate a novel backward type inference system for a large fragment of XQuery, including all the XPath axis expressions. In particular, we show that our backward type inference for XPath axes is exact and its complexity is simple exponential.
- We prove soundness of our backward type inference system for the XQuery core, from which we can obtain a typechecking algorithm. We also formally analyze the complexity of our inference system, and show that its complexity is double exponential in terms of the given expression.

## 2. Syntax and semantics of an XQuery core

In this section, we introduce an XQuery core, a minimal XQuery fragment supporting all the navigational XPath axes. Our XQuery core is an extension of miniXQuery proposed in [8] with non-downward axes.

### 2.1. Focused trees

We first define XML trees as focused trees, inspired by Huet’s zipper data structure [18]. A focused tree is an XML node with its context: the siblings and the parent of the node, including the parent’s context recursively. Intuitively a context

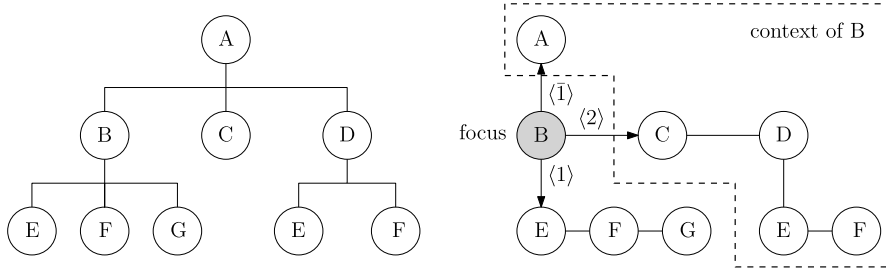


Fig. 1. An example XML tree structure and its corresponding binary representation.

records the path covered when traversing an XML tree from its root to a certain node. Thus focused trees allow us to easily navigate XML trees in any direction: both forward and backward navigation.

Below we formally define the syntax of our data model. We assume an alphabet  $\Sigma$  of labels, ranged over by  $\sigma$ .

Trees	$t ::= \sigma[tl]$
Tree lists	$tl ::= \epsilon \mid t :: tl$
Contexts	$c ::= \text{Top} \mid (tl; c[\sigma]; tl)$
Focused trees	$f ::= (t, c)$

A focused tree  $(t, c)$  is a pair consisting of a focused node (or a tree)  $t$  and its context  $c$ . A context  $c$  is Top if the focused node is at the root. Otherwise it is a triple  $(tl_l; c[\sigma]; tl_r)$ :  $tl_l$  is a list of the left siblings of the current focused node in reverse order (the first element of the list is the tree immediately to the left of the current node),  $c[\sigma]$  the context above the current node where  $\sigma$  is the label of the parent, and  $tl_r$  a list of the right siblings.

We now describe how to navigate a focused tree in a binary fashion. Given a focused tree  $f$ , forward navigation  $f \langle 1 \rangle$  and  $f \langle 2 \rangle$  respectively change the focus to the leftmost child and to the next right sibling of the current focused node. Conversely backward navigation  $f \langle \bar{1} \rangle$  and  $f \langle \bar{2} \rangle$  respectively change the focus to the parent and the preceding left sibling of the current node. In particular,  $f \langle \bar{1} \rangle$  is defined if and only if the current node is the leftmost node, i.e., it has no left sibling. Definition 2.1 formally defines the navigation of focused trees.

**Definition 2.1** (Navigation of focused trees).

$$\begin{aligned}
 (\sigma[t :: tl], c) \langle 1 \rangle &\stackrel{\text{def}}{=} (t, (\epsilon; c[\sigma]; tl)) \\
 (t, (tl_l; c[\sigma]; t' :: tl_r)) \langle 2 \rangle &\stackrel{\text{def}}{=} (t', (t :: tl_l; c[\sigma]; tl_r)) \\
 (t, (\epsilon; c[\sigma]; tl)) \langle \bar{1} \rangle &\stackrel{\text{def}}{=} (\sigma[t :: tl], c) \\
 (t', (t :: tl_l; c[\sigma]; tl_r)) \langle \bar{2} \rangle &\stackrel{\text{def}}{=} (t, (tl_l; c[\sigma]; t' :: tl_r))
 \end{aligned}$$

If the focused tree does not have the required shape, these operations are undefined.

**Example 2.2.** Consider the example XML tree in Fig. 1. If the node labeled  $B$  is a focus, then the focused tree  $f_B$  and its navigation is defined as follows. Below, for simplicity, for each leaf node, we write only its label. For instance, we write  $E$  instead of  $E[\epsilon]$ .

$$\begin{aligned}
 f_B &= (B[E :: F :: G :: \epsilon], (\epsilon; \text{Top}[A]; C :: D[E :: F :: \epsilon] :: \epsilon)) \\
 f_B \langle 1 \rangle &= (E, (\epsilon; c[B]; F :: G :: \epsilon)) \quad \text{where } c = (\epsilon; \text{Top}[A]; C :: D[E :: F :: \epsilon] :: \epsilon) \\
 f_C = f_B \langle 2 \rangle &= (C, (B[E :: F :: G :: \epsilon] :: \epsilon; \text{Top}[A]; D[E :: F :: \epsilon] :: \epsilon)) \\
 f_B \langle \bar{1} \rangle &= (A[B[E :: F :: G :: \epsilon] :: C :: D[E :: F :: \epsilon] :: \epsilon], \text{Top}) \\
 f_B \langle \bar{2} \rangle &= \text{undefined} \\
 f_D = f_C \langle 2 \rangle = f_B \langle 2 \rangle \langle 2 \rangle &= (D[E :: F :: \epsilon], (C :: B[E :: F :: G :: \epsilon] :: \epsilon; \text{Top}[A]; \epsilon))
 \end{aligned}$$

Note that the focused trees  $f_C$  and  $f_D$  focusing on the nodes labeled  $C$  and  $D$ , respectively, can be obtained by navigating  $f_B$  in a forward direction.

## 2.2. XQuery core

Fig. 2 defines the abstract syntax of a simplified navigational fragment of the XQuery core. In the XQuery core, which is defined in the XQuery 1.0 and XPath 2.0 Formal Semantics [3], navigational (i.e., structural) expressions are well separated from data value expressions (e.g., ordering and node identity testing) which make typechecking undecidable (see for

Expressions	$e ::= \epsilon \mid \langle \sigma \rangle \{e\} \langle /\sigma \rangle : u \mid e, e \mid \text{for } \$v \text{ in } e \text{ return } e$ $\mid \text{let } \$\bar{v} := e \text{ return } e \mid \text{if } \text{empty}(e) \text{ then } e \text{ else } e$ $\mid \$v/\text{axis}::n \mid \$var$
Variables	$\$var ::= \$v \mid \$\bar{v} \mid \$doc$
Axis names	$\text{axis} ::= \text{self} \mid \text{child} \mid \text{desc} \mid \text{fsibl} \mid \text{parent} \mid \text{anc} \mid \text{psibl}$
Name tests	$n ::= \sigma \mid *$
Values	$s ::= \epsilon \mid f :: s$

Fig. 2. Syntax of a navigational fragment of the XQuery core.

instance [23]). Since the full language of XQuery can be compiled into the XQuery core [3] and we are mainly interested in typechecking, we consider only navigational expressions in this paper.

First of all, we assume that an XML element constructor  $\langle \sigma \rangle \{e\} \langle /\sigma \rangle$  is always annotated with a type  $u$  (the precise definition of  $u$  is given in Section 3.1). In other words,  $\langle \sigma \rangle \{e\} \langle /\sigma \rangle : u$  in Fig. 2 can be considered as a combination of XQuery's untyped element constructor and `validate` expressions. In XQuery, the result of a construction expression  $\langle \sigma \rangle \{e\} \langle /\sigma \rangle$  is considered untyped (both statically and dynamically) unless it is validated using a `validate` expression. The `validate` expression checks if the constructed XML element conforms to the expected type at runtime, and if not, it raises a dynamic type error. Our element constructor  $\langle \sigma \rangle \{e\} \langle /\sigma \rangle : u$  differs from XQuery's `validate` expression in that its typechecking is done not dynamically but statically. For XQuery's untyped element constructors, i.e., without `validate`, we simply assume that they are annotated with `AnyElt` which is the type of all XML elements.

As for other expressions,  $\$doc$  is a special variable for reading the input document, and  $\epsilon$  denotes an empty sequence, i.e.,  $e, \epsilon = \epsilon, e = e$ . In a for-loop expression, an item variable  $\$v$  is bound to a single element node (or a single “item” in the XQuery terminology), whereas in a let-binding expression, a sequence variable  $\$\bar{v}$  is bound to a possibly empty sequence of nodes. In a conditional expression `if empty( $e$ ) then  $e_1$  else  $e_2$` , if the condition  $e$  evaluates to a non-empty sequence of nodes, then  $e_1$  is evaluated; otherwise,  $e_2$  is evaluated. An axis expression  $\$v/\text{axis}::n$  extracts the nodes that are reachable from the current node  $\$v$  through  $\text{axis}$  and that also satisfy the name test  $n$ . Path navigation can start only from an item variable. A name test  $n$  is either a node label  $\sigma$  or a wildcard pattern  $*$  that matches any label. For path navigation, we consider only `self`, `child`, `desc`, `fsibl`, `parent`, `anc`, and `psibl` axes because other axes can easily be encoded. (We use abbreviated names instead of the full name of the XPath axes.) We use the following syntactic sugar:

$$\$v/\text{desc-or-self}::n \equiv \$v/\text{self}::n, \$v/\text{desc}::n$$

An XQuery expression  $e$  evaluates to a value  $s$ , which is defined as a sequence of focused trees. This definition of values allows us to define the semantics in a compositional way. We write  $[f_1, \dots, f_n]$  for  $f_1 :: \dots :: f_n :: \epsilon$  and  $s_1, s_2$  for a sequence concatenation of  $s_1$  and  $s_2$ . In XQuery, all values are sequences and a single item (or tree) is considered a singleton sequence that contains only that item (or tree). Hence in the rest of the paper we use  $f$  and  $[f]$  interchangeably.

### 2.3. Semantics

Fig. 3 shows the semantics of the XQuery core, which is defined using the following denotation function:

$$\llbracket \_ \rrbracket : \text{Expression} \rightarrow \text{Substitution} \rightarrow \text{Value}$$

where a substitution  $\eta$  is a mapping from variables to values.

While most of the rules are straightforward and compositional, we took special care for an element constructor  $\langle \sigma \rangle \{e\} \langle /\sigma \rangle : u$ . First, suppose that the inner expression  $e$  evaluates to a sequence  $[f_1, \dots, f_n]$  of focused trees, where  $f_i = (t_i, c_i)$ . Then, we embed them into a new tree structure, namely  $\sigma[f_1 :: \dots :: f_n :: \epsilon]$ , whose context is `Top`. When navigating it, we need to update the context with respect to the new tree node. Therefore, we remove the old context from each focused tree  $f_i$  and obtain  $f = (\sigma[t_1 :: \dots :: t_n :: \epsilon], \text{Top})$ .

To evaluate a for-loop expression `for  $\$v$  in  $e_1$  return  $e_2$`  with substitution  $\eta$ , we first evaluate  $\llbracket e_1 \rrbracket_\eta$ . If the result is not an empty sequence, say  $[f_1, \dots, f_n]$ , then for each focused tree  $f_i$ , we evaluate the for-loop body  $e_2$  with an extended substitution  $\eta, \$v \mapsto f_i$ . Finally, we concatenate the results of evaluating  $\llbracket e_2 \rrbracket_{\eta, \$v \mapsto f_i}$  for  $i = 1, \dots, n$  in order. In contrast, if  $\llbracket e_1 \rrbracket_\eta$  evaluates to an empty sequence, then the for-loop expression also evaluates to an empty sequence.

To evaluate an axis expression  $\$v/\text{axis}::n$ , we analyze the shape of the focused tree bound to the for-loop variable  $\$v$ . The definition of  $\llbracket f/\text{axis}::n \rrbracket$  follows from the intuition behind the axis  $\text{axis}$ . For example,  $\llbracket f/\text{self}::n \rrbracket$  evaluates to a singleton sequence  $[f]$  if and only if the label of  $f$  matches the name test  $n$ . The semantics of `child` is defined using `self` and `fsibl` applied to the left-most child node. Note that  $f(\bar{1})$  and  $f(\bar{2})$  are never both defined for the same  $f$  and thus the definitions for the semantics of `parent` are mutually exclusive (the same is true for `anc`).  $\llbracket f/\text{fsibl}::n \rrbracket$  and  $\llbracket f/\text{psibl}::n \rrbracket$  recursively apply `fsibl` and `psibl` to the following and preceding siblings of  $f$ , respectively, if there exists such a node.  $\llbracket f/\text{desc}::n \rrbracket$  applies `self` and `desc` recursively to each child node of  $f$  and concatenates the results into a sequence.

$\llbracket \epsilon \rrbracket_\eta = \epsilon$	
$\llbracket \langle \sigma \rangle \{e\} \langle / \sigma \rangle : u \rrbracket_\eta = (\sigma[t_1 :: \dots :: t_n :: \epsilon], \text{Top})$	if $\llbracket e \rrbracket_\eta = [(t_1, c_1), \dots, (t_n, c_n)]$
$\llbracket e_1, e_2 \rrbracket_\eta = \llbracket e_1 \rrbracket_\eta, \llbracket e_2 \rrbracket_\eta$	
$\llbracket \$var \rrbracket_\eta = \eta(\$var)$	
$\llbracket \text{for } \$v \text{ in } e_1 \text{ return } e_2 \rrbracket_\eta = \Pi_{f_1, \dots, f_n} \llbracket e_2 \rrbracket_{\eta, \$v \mapsto f_i}$	if $\llbracket e_1 \rrbracket_\eta = [f_1, \dots, f_n]$
$\llbracket \text{for } \$v \text{ in } e_1 \text{ return } e_2 \rrbracket_\eta = \epsilon$	if $\llbracket e_1 \rrbracket_\eta = \epsilon$
$\llbracket \text{let } \$\bar{v} := e_1 \text{ return } e_2 \rrbracket_\eta = \llbracket e_2 \rrbracket_{\eta, \$\bar{v} \mapsto \llbracket e_1 \rrbracket_\eta}$	
$\llbracket \text{if nempty}(e) \text{ then } e_1 \text{ else } e_2 \rrbracket_\eta = \llbracket e_1 \rrbracket_\eta$	if $\llbracket e \rrbracket_\eta = f, s$
$\llbracket \text{if nempty}(e) \text{ then } e_1 \text{ else } e_2 \rrbracket_\eta = \llbracket e_2 \rrbracket_\eta$	if $\llbracket e \rrbracket_\eta = \epsilon$
$\llbracket \$v/axis::n \rrbracket_\eta = \llbracket \eta(\$v)/axis::n \rrbracket$	
$\llbracket f/self::n \rrbracket = [f]$	if $\text{name}(f) = n$ or $n = *$
$\llbracket f/self::n \rrbracket = \epsilon$	if $\text{name}(f) \neq n$ and $n \neq *$
$\llbracket f/child::n \rrbracket = \llbracket f'/self::n \rrbracket, \llbracket f'/fsibl::n \rrbracket$	if $f' = f(1)$
$\llbracket f/child::n \rrbracket = \epsilon$	if $f = (\sigma[\epsilon], c)$
$\llbracket f/parent::n \rrbracket = \llbracket f'/self::n \rrbracket$	if $f' = f(\bar{1})$
$\llbracket f/parent::n \rrbracket = \llbracket f'/parent::n \rrbracket$	if $f' = f(\bar{2})$
$\llbracket f/parent::n \rrbracket = \epsilon$	if $f = (t, \text{Top})$
$\llbracket f/fsibl::n \rrbracket = \llbracket f'/self::n \rrbracket, \llbracket f'/fsibl::n \rrbracket$	if $f' = f(2)$
$\llbracket f/fsibl::n \rrbracket = \epsilon$	if $f = (t, (t, \sigma[c]; \epsilon))$
$\llbracket f/psibl::n \rrbracket = \llbracket f'/psibl::n \rrbracket, \llbracket f'/self::n \rrbracket$	if $f' = f(\bar{2})$
$\llbracket f/psibl::n \rrbracket = \epsilon$	if $f = (t, (\epsilon; \sigma[c]; tl))$
$\llbracket f/anc::n \rrbracket = \llbracket f'/anc::n \rrbracket, \llbracket f'/self::n \rrbracket$	if $f' = f(\bar{1})$
$\llbracket f/anc::n \rrbracket = \llbracket f'/anc::n \rrbracket$	if $f' = f(\bar{2})$
$\llbracket f/anc::n \rrbracket = \epsilon$	if $f = (t, \text{Top})$
$\llbracket f/desc::n \rrbracket = \Pi_{f_1, \dots, f_m} \llbracket f_i/self::n \rrbracket, \llbracket f_i/desc::n \rrbracket$	if $\llbracket f/child::* \rrbracket = [f_1, \dots, f_m]$
$\llbracket f/desc::n \rrbracket = \epsilon$	if $\llbracket f/child::* \rrbracket = \epsilon$

Auxiliary definitions:  $\text{name}((\sigma[tl], c)) = \sigma$

Fig. 3. Semantics of the XQuery core.

### 3. Type language

Our type language is based on regular tree types [2] and a tree logic, which is a sub-logic of the alternation free modal  $\mu$ -calculus with converse [20]. In this section, we first briefly introduce regular tree types and the tree logic, together with their semantics in terms of sets of focused trees. Then we introduce our type language, regular tree types enriched with tree logic formulas [19].

#### 3.1. Regular tree types

We use a slight variant of XDuce's regular expression type language [24] to type sequences of XML trees (or *elements*), which is expressive enough to capture standard XML types such as DTD and XML Schema [25]. Formally we define our regular tree types as follows.

**Definition 3.1** (*Regular tree types*).

Unit types	$u ::= \text{element } n \{ \tau \}$
Name tests	$n ::= \sigma \mid *$
Sequence types	$\tau ::= u \mid () \mid \tau, \tau \mid (\tau \mid \tau) \mid \tau^* \mid \chi$

A sequence type  $\tau$  is a regular expression over unit types, where a unit type  $u$ , or a “prime type” in the XQuery terminology, corresponds to an XML element. (In general,  $u$  may also include primitive types such as `Int` or `String`, but for simplicity, we consider only element types.) As usual, we use the following abbreviations:  $\tau^+ \equiv \tau, \tau^*$  and  $\tau^? \equiv () \mid \tau$ . (We use  $\equiv$  both for syntactic equivalence and syntactic sugar.)

While the Kleene star  $*$  operator supports horizontal recursion, we use a *type environment* and type variables to support vertical recursion. A type environment  $E$  is a finite mapping from type variables  $x$  to types  $\tau$ . For example, we assume that every  $E$  that we consider in this paper maps a type variable `AnyElt` into `element * {AnyElt*}`, which is the type of all elements. The variables bound in  $E$  may be defined in a mutually recursive way, but recursion must be guarded by an element type to ensure well-formedness of types, *i.e.*, contractiveness of recursive types [26]. We also assume that regular expressions defined by  $E$  are composed of mutually exclusive unit types and *1-unambiguous* [27], which is standard and comes from XML Schema.

As usual, the semantics of regular tree types is defined as sets of forests, *i.e.*, sets of sequences of trees, and the subtyping relation is semantically defined as the set inclusion relation.

**Definition 3.2.** Given a type environment  $E$ , the semantics of types is defined by the smallest function  $\llbracket \_ \rrbracket_E$  that satisfies the following set of equations:

$$\begin{aligned}
\llbracket x \rrbracket_E &= \llbracket E(x) \rrbracket_E & \llbracket \tau^0 \rrbracket_E &= \{\epsilon\} \\
\llbracket () \rrbracket_E &= \{\epsilon\} & \llbracket \tau^{n+1} \rrbracket_E &= \llbracket \tau, \tau^n \rrbracket_E \\
\llbracket \tau \mid \tau' \rrbracket_E &= \llbracket \tau \rrbracket_E \cup \llbracket \tau' \rrbracket_E & \llbracket \text{element } \sigma \{ \tau \} \rrbracket_E &= \{ \llbracket \sigma[tl] \rrbracket_E \mid tl \in \llbracket \tau \rrbracket_E \} \\
\llbracket \tau^* \rrbracket_E &= \bigcup_{n \in \mathbb{N}} \llbracket \tau^n \rrbracket_E & \llbracket \text{element } * \{ \tau \} \rrbracket_E &= \{ \llbracket \sigma[tl] \rrbracket_E \mid \sigma \in \Sigma \text{ and } tl \in \llbracket \tau \rrbracket_E \} \\
\llbracket \tau, \tau' \rrbracket_E &= \{ [t_1, \dots, t_n, t'_1, \dots, t'_m] \mid [t_1, \dots, t_n] \in \llbracket \tau \rrbracket_E \text{ and } [t'_1, \dots, t'_m] \in \llbracket \tau' \rrbracket_E \}
\end{aligned}$$

Then, a type  $\tau_1$  is a subtype of  $\tau_2$ , denoted by  $\tau_1 <: \tau_2$ , if and only if  $\llbracket \tau_1 \rrbracket_E \subseteq \llbracket \tau_2 \rrbracket_E$ .

In the following, we assume that  $E$  is always well-formed and contains bindings for all variable references appearing in the types, and write  $\llbracket \tau \rrbracket$  as a shorthand for  $\llbracket \tau \rrbracket_E$ . We also assume that references  $x$  are implicitly replaced with their bindings at top level, so that a type  $\tau$  is really a regular expression of unit types.

The regular tree type language we gave above is standard and used to define the static type system in the XQuery standard and its various improvements in the literature. In such a type system, an XQuery expression is associated with a regular tree type, and the notion of a value (i.e., a sequence of tree nodes) *matching* a type can be defined as follows when nodes are represented as focused trees.

**Definition 3.3.** The *focused-tree interpretation*  $\llbracket \tau \rrbracket^\uparrow$  of a type  $\tau$  is defined as the set:

$$\{ [(t_1, c_1) \dots (t_n, c_n)] \mid [t_1 \dots t_n] \in \llbracket \tau \rrbracket \}$$

A value  $s$  is said to *match* a type  $\tau$  if  $s \in \llbracket \tau \rrbracket^\uparrow$ .

**Example 3.4.** Consider the example XML tree in Fig. 1 again. The focused trees  $f_B$ ,  $f_C$ , and  $f_D$  focusing on the nodes labeled  $B$ ,  $C$ , and  $D$ , respectively, defined in Example 2.2, match the following regular tree types  $\tau_B$ ,  $\tau_C$ , and  $\tau_D$ .

$$\begin{aligned}
\tau_B &= \text{element } B \{ \text{element } E \{ () \}, \text{element } F \{ () \}, \text{element } G \{ () \} \} \\
\tau_C &= \text{element } C \{ () \} \\
\tau_D &= \text{element } D \{ \text{element } E \{ () \}, \text{element } F \{ () \} \}
\end{aligned}$$

They also match a more general type such as  $\text{AnyElt}$  or  $\tau_{Node}$  which is defined as follows:

$$\begin{aligned}
\tau_{Node} &= \text{element } B \{ x^* \} \mid \text{element } C \{ () \} \mid \text{element } D \{ x^* \} \\
x &= (\text{element } E \{ x^* \}, \text{element } F \{ () \}) \mid \text{element } G \{ x^* \} \mid ()
\end{aligned}$$

Note however that we cannot describe the context information using Definition 3.1.

As shown in the above definition and example, regular tree types denote sequences of trees, and their interpretation is lifted to sequences of focused trees by simply ignoring the context part. In other words, using regular tree types, the type system cannot properly address expressions that analyze the shape of the context of a given focused tree: given  $f$  of type  $\tau$ , we cannot deduce a precise type for  $f(\bar{1})$ ,  $f(\bar{2})$ , and  $f(\bar{2})$  because when  $f = (t, c)$ ,  $\tau$  only contains information about the subtree  $t$ , but those expressions require information about the context  $c$ .

More specifically, consider an expression `for $v in e return $v/psibl::*`. Let us consider forward type inference; reasoning with backward type inference is similar. Suppose that  $e$  is of type  $\tau_{Node}$  and reduces to  $f_D$ . Then, we need to compute  $f_D/\text{psibl}::*$ , which reduces to  $[f_B, f_C]$ . The type of this result, however, should be determined by analyzing  $\tau_{Node}$  only, without evaluating the given expression. Since  $\tau_{Node}$  does not contain any useful information about its preceding siblings, we cannot deduce a meaningful type for  $f_D/\text{psibl}::*$ , and thus for the entire for-loop expression. Therefore, every type system for XQuery built solely on the type language given in Definition 3.1 simply gives to this expression the most general type  $\text{AnyElt}^*$ . To describe contexts and type navigational expressions precisely, we propose to use a tree logic in the next section.

### 3.2. A tree logic

To describe sets, i.e., types, of focused trees rather than just sets of trees, we use a variant of the logic language defined in [20]. The tree logic, defined below, is expressive enough to support all XQuery types, and the satisfiability problem for a logical formula of size  $n$  can efficiently be decided with an optimal  $2^{O(n)}$  worst-case time complexity bound [28].

**Definition 3.5** (Logic formulas).

$$\varphi, \psi ::= \top \mid \sigma \mid \neg\sigma \mid \alpha \mid \neg\alpha \mid \varphi \vee \psi \mid \varphi \wedge \psi \mid \langle a \rangle \varphi \mid \neg \langle a \rangle \varphi \mid X \mid \mu(X_i = \varphi_i)_{i \in I} \text{ in } \psi$$

$a \in \{1, 2, \bar{1}, \bar{2}\}$  are called *programs*, corresponding to the four directions where trees can be navigated. A program is used in an existential formula  $\langle a \rangle \varphi$ , denoting the existence of a subtree at the direction of  $a$  that satisfies the subformula  $\varphi$ . Other formulas include the truth predicate  $\top$ , atomic propositions  $\sigma$  (denoting the label of

the focused tree), node identifiers (a.k.a. *nominals*)  $\alpha$ , disjunction and conjunction of formulas, and least n-ary fixed points. In particular, an n-ary fixed point  $\mu(X_i = \varphi_i)_{i \in \{1, \dots, n\}}$  in  $\psi$  represents possibly mutually recursive definitions let  $\text{rec } X_1 = \varphi_1$  and  $X_2 = \varphi_2$  and  $\dots$  and  $X_n = \varphi_n$  in  $\psi$  (written in OCaml syntax) where  $\varphi_i$  and  $\psi$  may contain  $X_j$  for any  $i, j \in \{1, \dots, n\}$ . We also use the following abbreviations:  $\perp$  to mean  $\neg\top$ ,  $[a]\varphi$  for  $\neg(a)\top \vee (a)\varphi$ , and  $\mu X.\varphi$  for  $\mu(X = \varphi)$  in  $\varphi$ . The universal modality  $[a]\varphi$  encodes that a subtree at the direction of  $a$  does not exist, or else it satisfies  $\varphi$ . In this work, we consider only cycle-free formulas (for example,  $\mu(X = \langle 1 \rangle (\varphi \vee \langle \bar{1} \rangle X))$  in  $X$  is not cycle free) and thus the logic is closed under negation [28].

The semantics of a logical formula is defined as the set of focused trees such that the formula is satisfied at the current node. We use the following interpretation function:

$$\langle\langle - \rangle\rangle : \text{Formula} \rightarrow \text{Substitution} \rightarrow \text{FocusedTreeSet}$$

where a substitution  $V$  is a finite map from recursion variables to sets of focused trees. In the definition below, we use  $\mathcal{F}$  to denote the set of all focused trees and  $\text{name}(f)$  to denote the label at the current node of  $f$ . Moreover, to support nominals, we extend the syntax of focused trees so as to annotate a focused tree  $f$  with a set  $L$  of nominals. We write  $f^L$  to mean  $(t^L, c) = (\sigma^L[t], c)$  when  $f = (t, c)$  and  $t = \sigma[t]$ . In this work, nominals are purely a semantic notion and used only for typing descendants (see Section 4.2.6). Hence, for simplicity, we omit  $L$  whenever it is irrelevant to the discussion at hand, which is almost always the case.

**Definition 3.6** (Interpretation of formulas).

$$\begin{aligned} \langle\langle \top \rangle\rangle_V &\stackrel{\text{def}}{=} \mathcal{F} & \langle\langle X \rangle\rangle_V &\stackrel{\text{def}}{=} V(X) \\ \langle\langle \sigma \rangle\rangle_V &\stackrel{\text{def}}{=} \{f \mid \text{name}(f) = \sigma\} & \langle\langle (a)\varphi \rangle\rangle_V &\stackrel{\text{def}}{=} \{f \langle \bar{a} \rangle \mid f \in \langle\langle \varphi \rangle\rangle_V\} \\ \langle\langle \neg\sigma \rangle\rangle_V &\stackrel{\text{def}}{=} \{f \mid \text{name}(f) \neq \sigma\} & \langle\langle \neg(a)\top \rangle\rangle_V &\stackrel{\text{def}}{=} \{f \mid f \langle a \rangle \text{ undefined}\} \\ \langle\langle \alpha \rangle\rangle_V &\stackrel{\text{def}}{=} \{f^L \mid \alpha \in L\} & \langle\langle \varphi \vee \psi \rangle\rangle_V &\stackrel{\text{def}}{=} \langle\langle \varphi \rangle\rangle_V \cup \langle\langle \psi \rangle\rangle_V \\ \langle\langle \neg\alpha \rangle\rangle_V &\stackrel{\text{def}}{=} \{f^L \mid \alpha \notin L\} & \langle\langle \varphi \wedge \psi \rangle\rangle_V &\stackrel{\text{def}}{=} \langle\langle \varphi \rangle\rangle_V \cap \langle\langle \psi \rangle\rangle_V \\ \langle\langle \mu(X_i = \varphi_i)_{i \in I} \text{ in } \psi \rangle\rangle_V &\stackrel{\text{def}}{=} & & \\ \text{let } \mathcal{S} = \{(T_i)_{i \in I} \in \mathcal{P}(\mathcal{F})^I \mid \forall j \in I, \langle\langle \varphi_j \rangle\rangle_{V[\bar{T}_i/\bar{X}_i]} \subseteq T_j\} \text{ in} & & & \\ \text{let } \forall i \in I, U_i = \bigcap_{(T_j) \in \mathcal{S}} T_j \text{ in } \langle\langle \psi \rangle\rangle_{V[\bar{U}_i/\bar{X}_i]} & & & \\ \text{where } V[\bar{T}_i/\bar{X}_i](X) &\stackrel{\text{def}}{=} \begin{cases} V(X) & \text{if } X \notin \{X_i\}_{i \in I} \\ T_i & \text{if } X = X_i \end{cases} \end{aligned}$$

In the rest of the paper, we consider only closed formulas and write  $\langle\langle \psi \rangle\rangle$  for  $\langle\langle \psi \rangle\rangle_\emptyset$ . We say that a focused tree  $f$  matches a formula  $\psi$  if  $f \in \langle\langle \psi \rangle\rangle$ .

**Example 3.7.** Consider the focused trees  $f_B$ ,  $f_C$ , and  $f_D$  given in Example 2.2.  $f_D$  matches a formula  $\psi_D$  where the underlined part describes the subtree rooted at  $D$  and the other part describes its context.

$$\psi_D = \underline{D \wedge \langle 1 \rangle (E \wedge \langle 2 \rangle F)} \wedge \langle \bar{2} \rangle (C \wedge \langle \bar{2} \rangle (B \wedge \langle 1 \rangle (E \wedge \langle 2 \rangle (F \wedge \langle 2 \rangle G)) \wedge \langle \bar{1} \rangle A))$$

From  $\psi_D$ , we can now infer formulas  $\psi_C$  and  $\psi_B$  for  $f_C$  and  $f_B$  which are the preceding siblings of  $f_D$ .

$$\begin{aligned} \psi_C &= \langle 2 \rangle \psi_D = \langle 2 \rangle (D \wedge \langle 1 \rangle (E \wedge \langle 2 \rangle F)) \wedge \underline{C} \wedge \langle \bar{2} \rangle (B \wedge \langle 1 \rangle (E \wedge \langle 2 \rangle (F \wedge \langle 2 \rangle G)) \wedge \langle \bar{1} \rangle A) \\ \psi_B &= \langle 2 \rangle \langle 2 \rangle \psi_D = \langle 2 \rangle (\langle 2 \rangle (D \wedge \langle 1 \rangle (E \wedge \langle 2 \rangle F)) \wedge C) \wedge \underline{B} \wedge \langle 1 \rangle (E \wedge \langle 2 \rangle (F \wedge \langle 2 \rangle G)) \wedge \langle \bar{1} \rangle A \end{aligned}$$

### 3.3. Formula-enriched sequence types

In order to type sequences of focused trees, which are values of our XQuery core, we simply enrich the type language in Definition 3.1 by associating a formula to each unit type. The enriched types, which we call formula types, are thus regular expressions of pairs of unit types and formulas, as defined below.

**Definition 3.8** (Formula types).

$$\rho ::= (\varphi, u) \mid () \mid \rho, \rho \mid (\rho \mid \rho) \mid \rho^+$$

A formula type  $(\varphi, u)$  describes a focused tree  $(t, c)$  where  $u$  describes only  $t$  while  $\varphi$  may describe both  $t$  and  $c$ . The interpretation of a pair  $(\varphi, u)$  is defined as a set of singleton sequences of focused trees which match both  $\varphi$  and  $u$ :

$$\llbracket (\varphi, u) \rrbracket = \{ \llbracket (t, c) \rrbracket \mid t \in \llbracket u \rrbracket \text{ and } (t, c) \in \llbracket \langle \varphi \rangle \rrbracket \}$$

From this, the semantics of formula types in terms of sets of sequences of focused trees is defined in the obvious manner. Then, the subtyping relation  $\rho_1 <: \rho_2$  is semantically defined as the set inclusion relation  $\llbracket \rho_1 \rrbracket \subseteq \llbracket \rho_2 \rrbracket$ .

**Example 3.9.** Consider a for-loop expression `for $v in $f return $v/psibl::*`. If the type of  $f_D$  is given as  $(\psi_D, \tau_D)$ , then the type of the whole expression may be deduced as  $(\mu X. (2) (\psi_D \vee X), \text{AnyElt})^*$ . The type states that the for-loop expression will reduce to a possibly empty sequence of focused trees, each of which is of type `AnyElt` and has a following sibling that satisfies  $\psi_D$ . As discussed in Section 3.1, we cannot deduce any meaningful information in the unit type part, and thus simply use `AnyElt`.

The rationale behind the use of formula types is that it provides more flexibility. From the above example, one might think that regular expressions of formulas would be sufficient, which is true for backward type inference for XPath axes. However, sometimes, we may want to ignore context information, for example, to construct a new XML tree node using existing focused trees. In this case, we need to eliminate the context information from the formula matched with each focused tree. Unfortunately, it is nontrivial to eliminate only context information in the presence of recursive formulas. Thus, by combining formulas with unit types, we can make use of the usual unit type part to eliminate the context information, ignoring the formula part, if necessary. Moreover, although we do not investigate in this paper, by using the same type language as in [19], it would be easier to integrate our backward type inference with their forward type inference.

In Section 2.2, we assumed that every XML element constructor was annotated not with a formula type  $(\varphi, u)$  but with a unit type  $u$ . The reason is that an element constructor always reduces to a single tree node whose context is `Top`, and thus there is no need to use a formula type for the annotation. We simply consider  $u$  to be  $(\top, u)$ .

#### 4. Inference for XPath axes

In this section, we present a sound and complete backward type inference system for XPath axes, and based on this we will develop a backward type inference system for the XQuery core in Section 5. In backward type inference, we are given an expression  $e$  and an output type  $\rho_o$  for a sequence of focused trees that  $e$  may produce. Then we infer an input type  $\rho_i$  such that for any focused tree  $f$ , the following conditions hold.

- (Soundness) If  $f$  is of type  $\rho_i$ , then  $e(f)$  produces a sequence of nodes of type  $\rho_o$ .
- (Completeness) If  $e(f)$  produces a sequence of nodes of type  $\rho_o$ , then  $f$  is of type  $\rho_i$ .

When considering XPath axes, from the soundness perspective, we infer a type describing a set of input trees such that when applied to an axis, each input tree produces a sequence of nodes that has the output type  $\rho_o$ . Moreover, since XPath axes can only be applied to a for-loop variable in our XQuery core, from the completeness perspective, we infer from a given axis  $axis$  and an output type  $\rho$ , a single formula type  $(\varphi, u)$  (possibly their union) that the input tree, *i.e.*, the for-loop variable, must satisfy. In particular, we design the inference rules in such a way that the following invariant holds.

**Invariant 4.1.** *In our backward type inference system for XPath axes, if  $(\varphi_1, u_1) \mid \dots \mid (\varphi_n, u_n)$  is an inferred input type, a subtype relation  $\varphi_i <: u_i$  holds for all  $i$ 's, that is, for any  $t$  and  $c$ , if  $(t, c) \in \llbracket \langle \varphi_i \rangle \rrbracket$ , then  $\llbracket t \rrbracket \in \llbracket u_i \rrbracket$ , or equivalently,  $\llbracket \langle \varphi_i \rangle \rrbracket \subseteq \llbracket u_i \rrbracket \uparrow$  by identifying a focused tree  $f$  with a singleton sequence  $\llbracket f \rrbracket$ .*

The implication of this invariant is that for type inference for XPath axes, we can safely ignore the unit type part because it is always less precise than the formula part. Still, the unit type part is useful when typing XQuery expressions such as element constructors and thus we infer a useful unit type for XPath axes whenever possible.

Formally, the subtype relation  $\varphi <: u$  between formula  $\varphi$  and unit type  $u$  can be checked in two steps. First, we translate  $u$  into a downward-only formula which is true at any tree node matching this unit type, regardless of its context. In other words, we translate  $u$  into a formula which holds at any node of an XML tree if and only if the tree rooted at that node satisfies  $u$ . Technically, this translation can be done using an auxiliary function  $\text{form}(u)$ , which is defined and proved correct in [19]. (For its precise definition, we refer the reader to Figure 10 in [19].) Next, we test the satisfiability of the formula  $\varphi \wedge \neg \text{form}(u)$ , for example, using the decision procedure presented in [28]; in fact,  $\llbracket \langle \varphi \wedge \neg \text{form}(u) \rangle \rrbracket = \emptyset$  if and only if any focused tree matching  $\varphi$  also satisfies  $u$ , *i.e.*,  $\llbracket \langle \varphi \rangle \rrbracket \subseteq \llbracket u \rrbracket \uparrow$ .

Below we present inference rules using a judgment of the form  $\rho_i \leftarrow axis : n, \rho_o$  where input type  $\rho_i$  is always of the form  $(\varphi_1, u_1) \mid \dots \mid (\varphi_n, u_n)$ . We first look into the inference rules for `self` and `parent`.

##### 4.1. Inference rules for `self` and `parent`

###### 4.1.1. Self

Fig. 4 shows inference rules for `self`. Basically `self::n` returns a singleton sequence containing the input tree if it satisfies the name test  $n$ ; otherwise it returns an empty sequence. Conversely, if the output type is  $()$ , it means that



<p style="text-align: center; margin: 0;"><b>SELF-EMPTY</b></p> $\frac{}{\neg k(n), \text{AnyElt} \leftarrow \text{self} :: n, ()}$ <p style="text-align: center; margin: 0;"><b>SELF-SEQ1</b></p> $\frac{\neg \text{nullable}(\rho_i) \quad \text{nullable}(\rho_j) \quad \rho' \leftarrow \text{self} :: n, \rho_i}{\rho' \leftarrow \text{self} :: n, (\rho_1, \rho_2)} \quad (i, j = 1, 2, i \neq j)$ <p style="text-align: center; margin: 0;"><b>SELF-SEQ2</b></p> $\frac{\text{nullable}(\rho_1) \quad \text{nullable}(\rho_2) \quad \rho'_i \leftarrow \text{self} :: n, \rho_i}{\rho'_1   \rho'_2 \leftarrow \text{self} :: n, (\rho_1, \rho_2)} \quad (i = 1, 2)$ <p style="text-align: center; margin: 0;"><b>SELF-PLUS</b></p> $\frac{\rho' \leftarrow \text{self} :: n, \rho}{\rho' \leftarrow \text{self} :: n, \rho^+}$	<p style="text-align: center; margin: 0;"><b>SELF-FORMULA</b></p> $\frac{}{(\varphi \wedge k(n) \wedge \text{form}(u), u) \leftarrow \text{self} :: n, (\varphi, u)}$ <p style="text-align: center; margin: 0;"><b>SELF-OR</b></p> $\frac{\rho'_i \leftarrow \text{self} :: n, \rho_i}{\rho'_1   \rho'_2 \leftarrow \text{self} :: n, (\rho_1   \rho_2)} \quad (i = 1, 2)$ <p style="text-align: center; margin: 0;"><b>PARENT</b></p> $\frac{\rho' \leftarrow \text{self} :: n, \rho}{\text{child-type}(\rho') \leftarrow \text{parent} :: n, \rho}$
--	---

**Auxiliary definitions:**

$k(*) = \top$ $\text{nullable}() = \text{true}$ $\text{nullable}(\varphi, u) = \text{false}$ $\text{nullable}(\rho^+) = \text{nullable}(\rho)$ $\text{Prime}() = ()$ $\text{Prime}(u) = u$ $\text{Prime}(\tau^+) = \text{Prime}(\tau)$ $\text{distrib}(\chi, ()) = ()$ $\text{distrib}(\chi, u) = (\chi \wedge \text{form}(u), u)$ $\text{distrib}(\chi, \tau^+) = \text{distrib}(\chi, \tau)^+$	$k(\sigma) = \sigma$ $\text{nullable}(\rho_1, \rho_2) = \text{nullable}(\rho_1) \wedge \text{nullable}(\rho_2)$ $\text{nullable}(\rho_1   \rho_2) = \text{nullable}(\rho_1) \vee \text{nullable}(\rho_2)$ $\text{Prime}(\tau_1, \tau_2) = \text{Prime}(\tau_1)   \text{Prime}(\tau_2)$ $\text{Prime}(\tau_1   \tau_2) = \text{Prime}(\tau_1)   \text{Prime}(\tau_2)$ $\text{distrib}(\chi, (\tau_1, \tau_2)) = (\text{distrib}(\chi, \tau_1), \text{distrib}(\chi, \tau_2))$ $\text{distrib}(\chi, \tau_1   \tau_2) = (\text{distrib}(\chi, \tau_1)   \text{distrib}(\chi, \tau_2))$ $\text{child-type}(\rho_1   \rho_2) = \text{child-type}(\rho_1)   \text{child-type}(\rho_2)$ $\text{child-type}(\neg k(n), \text{AnyElt}) = (\text{has-parent}(\neg k(n)) \vee \varphi_{\text{root}}, \text{AnyElt})$ $\text{child-type}((\varphi, \text{element } n \{ \tau \})) = \text{distrib}(\text{has-parent}(\varphi), \text{Prime}(\tau))$ $\text{has-parent}(\chi) = \mu Z. (\bar{1}) \chi \vee (\bar{2}) Z$ $\varphi_{\text{root}} = \neg(\bar{1}) \top \wedge \neg(\bar{2}) \top \wedge \neg(2) \top$
---	---

**Fig. 4.** Inference rules for self and parent.

the input tree fails the name test and thus has type  $\neg k(n)$  (rule SELF-EMPTY). Here  $k(n)$  is the translation of  $n$  into a corresponding formula, i.e.,  $k(*) = \top$  and  $k(\sigma) = \sigma$ .

If the output type is a single formula type  $(\varphi, u)$ , it means that the input tree has that type: more precisely, the input tree should satisfy both  $\varphi$  and  $k(n)$ , and at the same time should have type  $u$  (rule SELF-FORMULA). All these constraints are encoded in the formula  $\varphi \wedge k(n) \wedge \text{form}(u)$  where we translate the unit type  $u$  into a formula using the function  $\text{form}(u)$ . In the rule SELF-FORMULA, since  $\varphi \wedge k(n) \wedge \text{form}(u) \prec: u$  holds, i.e.,  $\langle\langle \varphi \wedge k(n) \wedge \text{form}(u) \wedge \neg \text{form}(u) \rangle\rangle = \emptyset$ , Invariant 4.1 holds. In addition, when  $u = \text{element } \sigma \{ \tau \}$ , the inferred input formula is unsatisfiable if  $n = \sigma'$  and  $\sigma \neq \sigma'$ . In other words, there exist no tree nodes that produce a tree of type  $\text{element } \sigma \{ \tau \}$  when applied to  $\text{self} :: \sigma'$ , because no tree nodes can have different labels at the same time.

If the output type is a sequence type  $(\rho_1, \rho_2)$ , at least one type needs to be nullable (i.e., the interpretation of the type includes an empty sequence  $\epsilon$ ) since  $\text{self} :: n$  returns at most one tree as output. The type of the input tree is then the type inferred from the non-nullable part of the output type (rule SELF-SEQ1). If both  $\rho_1$  and  $\rho_2$  are nullable, we take the union of the input types inferred from them (rule SELF-SEQ2). When the output type is a union type, the input tree may also have a union type of the two, each of which is inferred from one summand of the output type (rule SELF-OR). Lastly, if the output type is a plus type  $\rho^+$ , the input type should be inferred from  $\rho$  since  $\text{self} :: n$  returns at most one node (rule SELF-PLUS).

**Proof of Invariant 4.1 for self.** By induction on a derivation of  $\rho' \leftarrow \text{self} :: n, \rho$ .  $\square$

**4.1.2. Parent**

The intuition behind type inference for parent is simple. Given an output type  $\rho$ , it is the type of the parent of the input context node. Moreover, if we infer  $\rho'$  using the inference rules for self with  $\rho$ , then the parent node is also of type  $\rho'$ . In other words, the input node is a child of the node of type  $\rho'$ . Therefore, for the input node, we extract a child type from  $\rho'$  using an auxiliary function  $\text{child-type}()$ .

To illustrate, assume that the output type  $\rho$  is given as  $(A, \tau_A)^+$  where  $\tau_A = \text{element } A \{ \tau_B, \tau_C, \tau_D \}$  and  $\tau_B, \tau_C$ , and  $\tau_D$  are defined in Example 3.4. Note that  $\rho$  can be used as a type for the focused tree  $f_A$  rooted at the node labeled  $A$  in Fig. 1. By applying the inference rules for self  $:: n$ , we obtain  $\rho' = (\varphi_A, \tau_A)$  where  $\varphi_A = A \wedge k(n) \wedge \text{form}(\tau_A)$ . Note that  $\rho'$  is also a type for  $f_A$ . Suppose that given an input node  $f$ ,  $f/\text{parent} :: n$  reduces to  $f_A$ . This means that  $f$  must be one

of  $f_B$ ,  $f_C$ , and  $f_D$ , which are the child nodes of  $f_A$ . To deduce the type of  $f$ , first consider the formula part. Any child of  $f_A$  matches a formula  $\varphi_i = \text{has-parent}(\varphi_A) = \mu Z. \langle \bar{1} \rangle \varphi_A \vee \langle \bar{2} \rangle Z$  which simply states that the context node has a parent matching  $\varphi_A$ . For the regular expression type part, we can deduce from  $\tau_A$  that any child of  $f_A$  matches a regular tree type  $\tau_B \mid \tau_C \mid \tau_D$ . Finally, by distributing  $\varphi_i$  over  $\tau_B \mid \tau_C \mid \tau_D$  using an auxiliary function  $\text{distrib}()$ , we obtain an input type  $(\varphi_i \wedge \text{form}(\tau_B), \tau_B) \mid (\varphi_i \wedge \text{form}(\tau_C), \tau_C) \mid (\varphi_i \wedge \text{form}(\tau_D), \tau_D)$ .

In general, given an output type  $\rho$ , when we infer a parent type  $(\varphi, \text{element } n \{ \tau \})$  using the inference rules for `self`,  $\tau$  may be an arbitrary regular expression. Therefore, for the regular expression type part, we compute a child type using an auxiliary function  $\text{Prime}(\tau)$  [8] which extracts all unit types at the top level of  $\tau$  and constructs their disjunction. Moreover, if the output type  $\rho$  is nullable, then the input node may be a root. In other words, if  $(\neg k(n), \text{AnyElt})$  is inferred for the parent node using the inference rules for `self`, then the input context node satisfies the formula  $\text{has-parent}(\neg k(n)) \vee \varphi_{\text{root}}$ , where  $\varphi_{\text{root}}$  is defined as  $\neg \langle \bar{1} \rangle \top \wedge \neg \langle \bar{2} \rangle \top \wedge \neg \langle 2 \rangle \top$  and specifies that a given node is a root. That is, the input context node has a parent whose label is different from  $n$ , or else it does not have a parent. Note that we cannot specify the fact that the input node may be a root in the regular tree type part.

**Proof of Invariant 4.1 for parent.** Suppose  $(\varphi_1, u_1) \mid \dots \mid (\varphi_k, u_k) \leftarrow \text{self} : : n, \rho$ . Then  $\varphi_i <: u_i$  for all  $i$ 's. For each  $i$ , we need to show that if  $\text{child-type}((\varphi_i, u_i)) = (\varphi'_1, u'_1) \mid \dots \mid (\varphi'_j, u'_j)$ , then  $\varphi'_j <: u'_j$  for all  $j$ 's. If  $\varphi_i = \neg k(n)$  and  $u_i = \text{AnyElt}$ , then the proof is straightforward since  $u'_j = \text{AnyElt}$ . Thus, suppose  $u_i = \text{element } n_i \{ \tau_i \}$  and  $\text{Prime}(\tau_i) = u'_1 \mid \dots \mid u'_l$ . Then, for each  $j$ ,  $\varphi'_j = \text{has-parent}(\varphi_i) \wedge \text{form}(u'_j)$  and therefore  $\varphi'_j <: u'_j$ .  $\square$

#### 4.2. Inference rules for other axes

Unlike `self` and `parent`, for other axes, given an output type  $(\varphi, u)$ , we cannot specify the exact shape of the input tree in the unit type part of the inferred input type since the output unit type  $u$  does not contain information about the context. Hence, for other axes, we approximate the unit type part in the inferred input type. Still, we do not lose any precision since the formula part of the input type is exact. In other words, for type inference for XPath axes, we can safely ignore the unit type part of the inferred input type (Invariant 4.1). Nevertheless, we try to infer a more precise type than `AnyElt` for the unit type part if possible. More precisely, we simply infer `AnyElt` for `psibl`, `fsibl`, and `desc`, while inferring a more precise type for `child` and `anc`. As studied in [8,19], in forward type inference systems using only regular tree types as its type language, one can infer precise types only for `self`, `child`, and `desc`. In contrast, in our backward type inference system, we infer precise regular tree types only for `self`, `parent`, `child`, and `anc` (the formula part is still exact for all XPath axes).

Another important difference between `self` and `parent` and other axes is that while the former requires us to inspect only a single node in the input tree, the rest of the axes requires us to inspect a sequence of nodes reached by navigating the axis from the input node and combine the constraints for all these nodes. In order to combine a set of constraints on a sequence of nodes, we use an additional judgment of the form  $\varphi \leftarrow \text{axis} : : n, \rho$  with  $\psi$ , for which we develop inference rules in such a way that the following property holds (the proof is given in Lemmas 4.5 and 4.6 in Section 4.3).

**Proposition 4.2.** Suppose  $\varphi \leftarrow \text{axis} : : n, \rho$  with  $\psi$  and  $\llbracket f/\text{axis} : : n \rrbracket = f_1, \dots, f_n$  for some input node  $f$ .

- For backward axes, let  $f = f_{n+1}$ . Then  $f \in \langle\langle \varphi \rangle\rangle$  if and only if  $\exists 1 \leq i \leq n+1$  such that  $f_i \in \langle\langle \psi \rangle\rangle$  and  $f_1, \dots, f_n \in \llbracket \rho \rrbracket$ .
- For forward axes, let  $f = f_0$ . Then  $f \in \langle\langle \varphi \rangle\rangle$  if and only if  $\exists 0 \leq i \leq n$  such that  $f_i \in \langle\langle \psi \rangle\rangle$  and  $f_1, \dots, f_n \in \llbracket \rho \rrbracket$ .

Specifically, given a backward (resp. forward) axis, the judgment  $\varphi \leftarrow \text{axis} : : n, \rho$  with  $\psi$  means that an input context node  $f$  satisfies the inferred input formula  $\varphi$  if and only if  $f/\text{axis} : : n$  produces a sequence of nodes  $f_1, \dots, f_n$  such that there exists some node  $f_i$  satisfying  $\psi$  and the node sequence  $f_i, \dots, f_n$  (resp.  $f_1, \dots, f_i$ ) is of type  $\rho$ . In particular, for a backward (resp. forward) axis,  $f_n$  (resp.  $f_1$ ) is the closest node satisfying the name test  $n$  at the direction of  $\text{axis}$  from the context node  $f$ . Moreover,  $f/\text{axis} : : n$  produces an empty sequence if and only if  $f$  satisfying  $\varphi$  also satisfies  $\psi$  and  $\rho$  is nullable. Note that using this judgment, we only infer a formula. We infer a unit type for the input node using auxiliary functions.

To illustrate the meaning of the judgment  $\varphi \leftarrow \text{axis} : : n, \rho$  with  $\psi$ , consider an example input tree node represented as a binary tree in Fig. 5. Suppose that  $\text{axis}$  is `psibl` and  $C$  is the context node. Then, there exists some node  $A$  reachable by navigating `psibl` from  $C$  that satisfies both the with parameter  $\psi$  and the name test  $n$ . Moreover, by analyzing  $\psi$ , we can obtain the constraints on the nodes reached by further navigating `psibl` from  $A$ . That is,  $\psi$  should contain the context information for  $A$  as its subformulas, for example, information on  $A$ 's preceding siblings. In addition, the sequence of the nodes from  $A$  to  $B$  that consists only of the nodes satisfying the name test  $n$  has type  $\rho$ , where  $B$  is the rightmost preceding sibling of  $C$  in document order (i.e., pre-order). Lastly, the context node  $C$  has the inferred type  $\varphi$ . In the subsection below, we give a more precise interpretation of the judgment when the output type is a sequence type of the form  $(\rho_1, \rho_2)$ .

With this interpretation, given  $\text{axis} : : n$  and an output type  $\rho$ , we first infer a formula  $\varphi$  using the judgment  $\varphi \leftarrow \text{axis} : : n, \rho$  with  $\psi_{\text{init}}$ , where the initial formula  $\psi_{\text{init}}$  is determined by  $\text{axis}$ . For example, if  $\text{axis}$  is `psibl`, then  $\psi_{\text{init}}$  is set to  $\mu X. [\bar{2}] (\neg k(n) \wedge X)$  which means that there exist no preceding siblings satisfying the name test  $n$ . Then, we compute

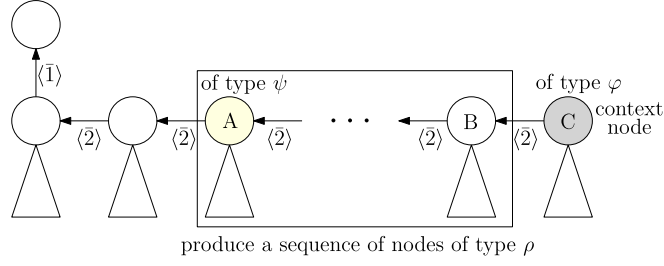


Fig. 5. Interpretation of  $\varphi \leftarrow \text{psibl} : : n, \rho$  with  $\psi$  when the context node is C: A and B are some nodes reached by navigating `psibl` from C.

$$\frac{\text{PSIBL} \quad \varphi \leftarrow \text{psibl} : : n, \rho \text{ with } \mu X. [\bar{2}] (\neg k(n) \wedge X)}{(\varphi, \text{AnyElt}) \leftarrow \text{psibl} : : n, \rho} \quad \frac{\text{PSIBL-FORMULA} \quad \varphi' = [\bar{2}] (\mu X. (\varphi \wedge k(n) \wedge \text{form}(u) \wedge \psi) \vee (\neg k(n) \wedge [\bar{2}] X))}{\varphi' \leftarrow \text{psibl} : : n, (\varphi, u) \text{ with } \psi}$$

Fig. 6. Inference rules for `psibl`.

$$\begin{array}{l} \text{AXIS-EMPTY} \\ \frac{}{\psi \leftarrow \text{axis} : : n, () \text{ with } \psi} \\ \\ \text{AXIS-OR} \\ \frac{\varphi_1 \leftarrow \text{axis} : : n, \rho_1 \text{ with } \psi}{\varphi_1 \vee \varphi_2 \leftarrow \text{axis} : : n, (\rho_1 | \rho_2) \text{ with } \psi} \quad (i = 1, 2) \\ \\ \text{AXIS-BACKWARD-SEQ} \\ \frac{\varphi_1 \leftarrow \text{axis} : : n, \rho_1 \text{ with } \psi \quad \varphi_2 \leftarrow \text{axis} : : n, \rho_2 \text{ with } \varphi_1}{\varphi_2 \leftarrow \text{axis} : : n, (\rho_1, \rho_2) \text{ with } \psi} \quad (\text{axis is psibl or anc}) \\ \\ \text{AXIS-FORWARD-SEQ} \\ \frac{\varphi_2 \leftarrow \text{axis} : : n, \rho_2 \text{ with } \psi \quad \varphi_1 \leftarrow \text{axis} : : n, \rho_1 \text{ with } \varphi_2}{\varphi_1 \leftarrow \text{axis} : : n, (\rho_1, \rho_2) \text{ with } \psi} \quad (\text{axis is fsibl or desc}) \\ \\ \text{AXIS-PLUS} \\ \frac{\varphi \leftarrow \text{axis} : : n, \rho \text{ with } X \vee \psi}{\mu X. \varphi \leftarrow \text{axis} : : n, \rho^+ \text{ with } \psi} \quad (X \text{ fresh}) \end{array}$$

Fig. 7. Common inference rules for `psibl`, `anc`, `fsibl`, and `desc`.

a unit type  $u$  using an appropriate auxiliary function depending on  $\text{axis}$ . When computing  $u$ , we ensure a subtype relation  $\varphi <: u$ . Finally, the input type is determined as a pair  $(\varphi, u)$  as shown below:

$$\frac{\varphi \leftarrow \text{axis} : : n, \rho \text{ with } \psi_{\text{init}} \quad u = \text{aux\_func}(\rho)}{(\varphi, u) \leftarrow \text{axis} : : n, \rho}$$

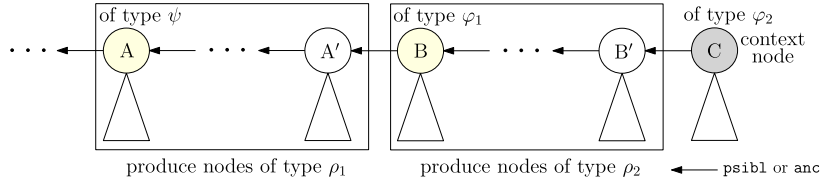
#### 4.2.1. Preceding siblings and generic inference rules

`psibl : : n` returns in document order the preceding siblings of the context node, say  $f$ , that satisfy the name test  $n$ . In other words, given an output type  $\rho$ , it denotes the type of the preceding siblings of  $f$ . Thus, the inferred formula for  $f$  obtained by analyzing  $\rho$  should accumulate the constraints (i.e., types) on its preceding siblings. The rest of the nodes reachable from  $f$  may have an arbitrary structure if they are not described by (the context part of) the output type  $\rho$ .

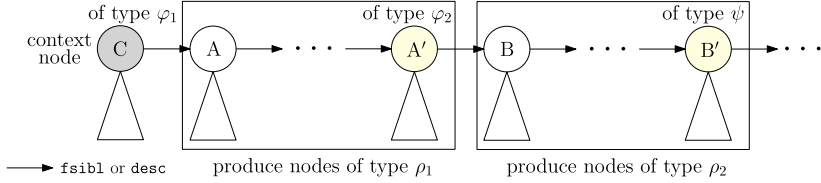
Fig. 6 shows the inference rules for `psibl`. In the rule `PSIBL`, we initially assume that there are no preceding siblings satisfying the name test  $n$ , that is,  $\mu X. [\bar{2}] (\neg k(n) \wedge X)$ . Then, we analyze the output type  $\rho$  using the judgment of the form  $\varphi \leftarrow \text{psibl} : : n, \rho$  with  $\psi$ . In this judgment,  $\psi$  is true at the leftmost preceding sibling returned by `psibl : : n` when the output type is  $\rho$  (e.g., the node A in Fig. 5). When the inferred input formula is  $\varphi$ , the final input type is a pair  $(\varphi, \text{AnyElt})$ . Since we cannot extract any meaningful information about the context node from the regular tree types of its preceding siblings, we simply use `AnyElt`. Thus, for `psibl`, Invariant 4.1 trivially holds.

When the output type is a single formula type  $(\varphi, u)$  and the with parameter is  $\psi$ , it means that there should be a preceding sibling satisfying the name test  $n$  such that both  $\varphi$  and  $\psi$  are true. Moreover, that sibling node should also have type  $u$ . All these constraints are encoded in the inferred formula  $\varphi'$  in the rule `PSIBL-FORMULA`. As in the rule `SELF-FORMULA`, we use function `form(u)` to translate the unit type  $u$  to a corresponding formula. In addition, since the initial with parameter given in the rule `PSIBL` guarantees that there are no preceding siblings satisfying the name test  $n$ , the two rules guarantee that if `psibl : : n` returns a single node, then the context node has only one preceding sibling satisfying  $n$ .

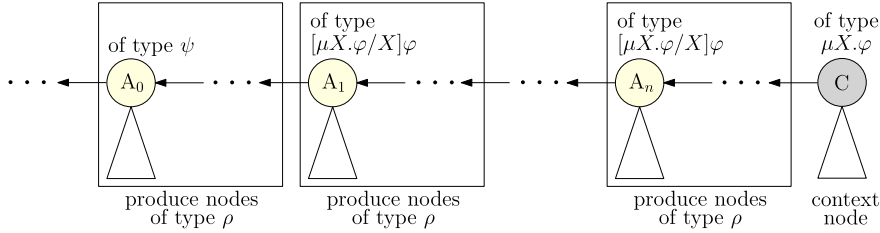
The rest of the inference rules for empty, sequence, union, and repetition types are generic and are also used for other axes—`anc`, `fsibl`, and `desc`. (When the output type is a sequence type, we distinguish backward axes from forward axes, and thus present two inference rules.) The common rules are given in Fig. 7. The first two rules are easy. If the output type is an empty type, the inferred input type is simply the formula  $\psi$  given as the with parameter (rule `AXIS-EMPTY`). Therefore, in combination with the rule `PSIBL`, the inferred formula in the rule `AXIS-EMPTY` specifies that no preceding sibling of the



**Fig. 8.** Interpretation of  $\varphi \leftarrow \text{axis} : : n, (\rho_1, \rho_2)$  with  $\psi$  when *axis* is a backward axis *psibl* or *anc*: we analyze the sequence type from left to right. We first infer  $\varphi_1$  and then  $\varphi_2$ .



**Fig. 9.** Interpretation of  $\varphi \leftarrow \text{axis} : : n, (\rho_1, \rho_2)$  with  $\psi$  when *axis* is a forward axis *fsibl* or *desc*: we analyze the sequence type from right to left. We first infer  $\varphi_2$  and then  $\varphi_1$ .



**Fig. 10.** Interpretation of  $\varphi \leftarrow \text{axis} : : n, \rho^+$  with  $\psi$  when *axis* is a backward axis. A similar illustration can be applied to forward axes.

input node should satisfy the name test. If the output type is a union type of two, we infer a formula from each and return the union of the two inferred formulas (rule *AXIS-OR*).

When the output type is a sequence type  $(\rho_1, \rho_2)$ , our analysis begins with the farthest node from the input context node among the nodes reached by navigating the given axis and proceeds towards the context node. Therefore, if the given axis is a backward axis such as *psibl* and *anc*, we analyze the output type from left to right (rule *AXIS-BACKWARD-SEQ*). More precisely, as depicted in Fig. 8, given a judgment  $\varphi \leftarrow \text{axis} : : n, (\rho_1, \rho_2)$  with  $\psi$ , we can conceptually divide the nodes reached by navigating *axis* from the context node C into two parts: the nodes from A to A' and those from B to B' that produce a sequence of nodes of type  $\rho_1$  and  $\rho_2$ , respectively, where the first part precedes the second part in document order. In particular,  $\psi$  is true at node A which is the first node in the first part. We first infer a formula  $\varphi_1$  from  $\rho_1$  and  $\psi$  using the judgment  $\varphi_1 \leftarrow \text{axis} : : n, \rho_1$  with  $\psi$ . Then  $\varphi_1$  is true at node B which is next to A' in document order and also the first node in the second part. Next, we infer a formula  $\varphi_2$  from  $\rho_2$  and  $\varphi_1$  using the judgment  $\varphi_2 \leftarrow \text{psibl} : : n, \rho_2$  with  $\varphi_1$ . Finally,  $\varphi_2$  is true at the context node and is returned as the input type.

The interpretation of the judgment  $\varphi \leftarrow \text{axis} : : n, (\rho_1, \rho_2)$  with  $\psi$  is dual if *axis* is a forward axis such as *fsibl* and *desc*. In this case, we analyze the output type from right to left, i.e.,  $\rho_2$  first (rule *AXIS-FORWARD-SEQ*). For example, as depicted in Fig. 9, with  $\rho_2$  and  $\psi$ , we start from a constraint on the last node B' at which  $\psi$  is true, among the nodes reached by navigating *axis* from the context node C, and subsequently infer constraints on the nodes appearing before B' in reverse order, e.g., from B through A' to A, until finally inferring the constraint on the context node.

When the output type is a repetition type  $\rho^+$ , we introduce a fresh recursion variable  $X$  (rule *AXIS-PLUS*). Then, we infer a formula  $\varphi$  from the output type  $\rho$  and the with parameter  $X \vee \psi$  using the judgment  $\varphi \leftarrow \text{axis} : : n, \rho$  with  $X \vee \psi$ . More precisely, as depicted in Fig. 10, there exists a block of nodes reached by navigating *axis* from the context node C, e.g., the nodes from  $A_n$  to the node before C, that produce a sequence of nodes of type  $\rho$ , each of which satisfies the name test  $n$ . Moreover, the lastly reached node  $A_n$  should satisfy  $X \vee \psi$ , while C should satisfy the inferred formula  $\varphi$  (where  $\varphi$  contains  $X \vee \psi$  as a subformula, for example, see the rule *PSIBL-FORMULA*). If  $A_n$  satisfied  $X$ , that is,  $[\mu X. \varphi / X] \varphi$ , there would be more blocks of nodes reached by further navigating *axis* from  $A_n$  that would produce nodes of type  $\rho$ , where the lastly reached node in each block, e.g.,  $A_1$ , would also satisfy  $X$ . This recursion terminates when some node satisfies  $\psi$  rather than  $X$ , e.g.,  $A_0$  (where the block of nodes containing  $A_0$  should also produce a sequence of nodes of type  $\rho$ ). Lastly, the closed recursive formula  $\mu X. \varphi$  is returned as the input type of the context node C.

ANC	ANC-FORMULA
$\varphi \leftarrow \text{anc} :: n, \rho \text{ with } \neg \text{has-anc}(k(n))$	$\varphi' = \mu X. \langle \bar{1} \rangle ((\varphi \wedge k(n) \wedge \text{form}(u) \wedge \psi) \vee (\neg k(n) \wedge X)) \vee \langle \bar{2} \rangle X$
$\text{distrib}(\varphi, \text{desc-type}(\rho)) \leftarrow \text{anc} :: n, \rho$	$\varphi' \leftarrow \text{anc} :: n, (\varphi, u) \text{ with } \psi$
$\text{has-anc}(\chi) =$	$\mu Z. \langle \bar{1} \rangle (\chi \vee Z) \vee \langle \bar{2} \rangle Z$
$\text{desc-type}(\cdot) =$	$()$
$\text{desc-type}(\langle \varphi, \text{element } n \{ \tau \} \rangle) =$	$\text{desc-type}(\tau)$
$\text{desc-type}(\rho_1, \rho_2) =$	$\text{desc-type}(\rho_1) \mid \text{desc-type}(\rho_2)$
$\text{desc-type}(\rho_1 \mid \rho_2) =$	$\text{desc-type}(\rho_1) \mid \text{desc-type}(\rho_2)$
$\text{desc-type}(\rho^+) =$	$\text{desc-type}(\rho)$
$\text{desc-type}(\text{element } n \{ \tau \}) =$	$\text{element } n \{ \tau \} \mid \text{desc-type}(\tau)$
$\text{desc-type}(\tau_1, \tau_2) =$	$\text{desc-type}(\tau_1) \mid \text{desc-type}(\tau_2)$
$\text{desc-type}(\tau_1 \mid \tau_2) =$	$\text{desc-type}(\tau_1) \mid \text{desc-type}(\tau_2)$
$\text{desc-type}(\tau^+) =$	$\text{desc-type}(\tau)$

Fig. 11. Inference rules for anc.

#### 4.2.2. Ancestors

Inference rules for  $\text{anc} :: n$  are the same as those for  $\text{psibl} :: n$  with two exceptions (they are both backward axes and use the same set of rules in Fig. 7): first the interpretation of the judgment and the initial value of the with parameter, and second the input type inferred when the output type is a single formula type  $(\varphi, u)$ . We briefly explain them in turn.

First, the interpretation of a judgment  $\varphi \leftarrow \text{anc} :: n, \rho \text{ with } \psi$  is as follows: there is a block of nodes reached by navigating  $\text{anc}$  from the context node such that it produces a sequence of nodes of type  $\rho$ , each of which satisfies the name test  $n$ . Moreover,  $\psi$  is true at the lastly reached node, or equivalently, the first node in document order, in that block. (We may reuse the example in Fig. 5 for  $\text{anc}$  by interpreting the left arrow in the figure as  $\langle \bar{1} \rangle$  followed by a possibly empty sequence of  $\langle \bar{2} \rangle$ s.) In the rule ANC in Fig. 11, we thus set the with parameter to  $\neg \text{has-anc}(k(n))$  to mean that there are no (more) ancestors satisfying the name test  $n$ .  $\text{has-anc}(\chi)$  is a formula that describes any tree node such that it has at least one ancestor at which  $\chi$  is true and  $\neg \text{has-anc}(\chi)$  is its negation.<sup>1</sup> Note that  $\langle \bar{2} \rangle$  denotes the left sibling of the context node if any, and  $\langle \bar{1} \rangle$  its parent if the context node has no left sibling and is not a root.

When the output type is  $(\varphi, u)$  and the with parameter is  $\psi$ , it means that the context node has an ancestor  $f$  that satisfies the name test  $n$  and is of type  $(\varphi, u)$  (rule ANC-FORMULA). Moreover,  $f$  should also satisfy  $\psi$ , which should contain as its subformulas the context information on the structure of  $f$ 's ancestors. The inferred input formula  $\varphi'$  is thus a recursive formula that denotes a tree node having an ancestor  $f$  satisfying all these constraints, i.e.,  $\varphi \wedge k(n) \wedge \text{form}(u) \wedge \psi$ . Furthermore, the ancestors between  $f$  and the context node should not satisfy the name test  $n$  and thus have type  $\neg k(n) \wedge X$ , which is also encoded in the inferred input formula  $\varphi'$ . Combined with the rule ANC, the rule ANC-FORMULA specifies that if the output type is a single formula type, then the context node has only one ancestor satisfying  $n$ .

As for the regular tree type part, we use an auxiliary function  $\text{desc-type}()$  which is a recursive version of  $\text{Prime}()$  and computes the type of all possible descendants. Note that for  $\text{anc}$ , the output type is the type of the ancestors of the input context node. In other words, the context node is one of their descendants. Hence, in the rule ANC, we distribute the inferred input formula over the union of all possible descendant types using the function  $\text{distrib}()$  defined in Fig. 4.

**Proof of Invariant 4.1 for anc.** In the rule ANC, suppose  $\text{distrib}(\varphi, \text{desc-type}(\rho)) = (\varphi_1, u_1) \mid \dots \mid (\varphi_k, u_k)$  for some  $k$ . We need to show  $\varphi_i <: u_i$  for all  $i$ 's. This is easy because  $\varphi_i = \varphi \wedge \text{form}(u_i)$  by the definition of  $\text{distrib}()$ .  $\square$

#### 4.2.3. Following siblings

$\text{fsibl}$  is the converse of  $\text{psibl}$ . To obtain the inference rules for  $\text{fsibl}$  in Fig. 12, we just replace  $\langle \bar{2} \rangle$  and  $[\bar{2}]$  in the rules PSIBL and PSIBL-FORMULA with  $\langle 2 \rangle$  and  $[2]$ , respectively, and use the rule AXIS-FORWARD-SEQ instead of the rule AXIS-BACKWARD-SEQ. More precisely, the interpretation of a judgment  $\varphi \leftarrow \text{fsibl} :: n, \rho \text{ with } \psi$  is as follows: there is an initial subsequence of the nodes reached by navigating  $\text{fsibl}$  from the context node that has type  $\rho$ , each of which satisfies the name test  $n$ . Moreover, the lastly reached node in that subsequence and the context node satisfy  $\psi$  and  $\varphi$ , respectively. Since our analysis always starts with the lastly reached node, i.e., the rightmost sibling in the case of  $\text{fsibl}$ , in the rule FSIBL, we set the initial with parameter to  $\mu X. [2] (\neg k(n) \wedge X)$  which means that there are no (more) following siblings satisfying the name test  $n$ . For the regular tree type part, we simply use AnyElt because of the lack of information about the context in the regular tree types of the following sibling nodes. Thus, for  $\text{fsibl}$ , Invariant 4.1 trivially holds.

If the output type is a single formula type  $(\varphi, u)$  and the with parameter is  $\psi$ , it means that one of the following siblings of the context node, say  $f$ , satisfies the name test  $n$  and is of type  $(\varphi, u)$  (rule FSIBL-FORMULA). Moreover,  $\psi$  should also be true at  $f$ . As with other axes,  $\psi$  should contain as its subformulas the context information on the structure of  $f$ 's following siblings. In addition, the following siblings between the context node and  $f$ , if any, should not satisfy the name test  $n$  and thus have type  $\neg k(n) \wedge \langle 2 \rangle X$ . All these constraints are encoded in the inferred input formula  $\varphi'$ . Note that in combination

<sup>1</sup> Technically this encoding allows the presence of hedges satisfying the formula (we do not impose the invariant that there is only a single root), but our semantics ensures that a formula accepts trees only.

$$\begin{array}{c}
\text{FSIBL} \\
\frac{\varphi \leftarrow \text{fsibl}::n, \rho \text{ with } \mu X.[2](\neg k(n) \wedge X)}{(\varphi, \text{AnyElt}) \leftarrow \text{fsibl}::n, \rho}
\end{array}
\qquad
\begin{array}{c}
\text{FSIBL-FORMULA} \\
\frac{\varphi' = (2)(\mu X.(\varphi \wedge k(n) \wedge \text{form}(u) \wedge \psi) \vee (\neg k(n) \wedge (2) X))}{\varphi' \leftarrow \text{fsibl}::n, (\varphi, u) \text{ with } \psi}
\end{array}$$

Fig. 12. Inference rules for fsibl.

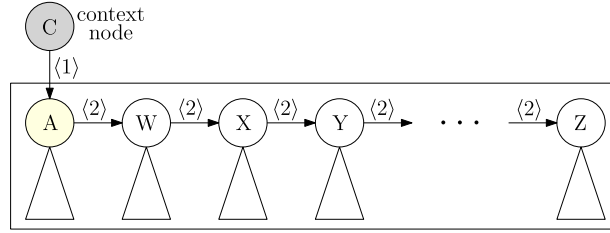


Fig. 13. child defined in terms of self-fsibl: the nodes in the box are C's children and A is its leftmost child. Therefore,  $\llbracket C/\text{child}::n \rrbracket = \llbracket A/\text{self-fsibl}::n \rrbracket$ . If  $\varphi \leftarrow \text{self-fsibl}::n, \rho$  and  $\varphi$  is true at node A, then  $\langle 1 \rangle \varphi$  is true at node C.

$$\begin{array}{c}
\text{CHILD-NULABLE} \\
\frac{\varphi \leftarrow \text{self-fsibl}::n, \rho \quad \text{nullable}(\rho)}{(\langle 1 \rangle \varphi, \text{parent-type}(\rho)) \leftarrow \text{child}::n, \rho}
\end{array}
\qquad
\begin{array}{c}
\text{CHILD-NOTNULL} \\
\frac{\varphi \leftarrow \text{self-fsibl}::n, \rho \quad \neg \text{nullable}(\rho)}{(\langle 1 \rangle \varphi, \text{parent-type}(\rho)) \leftarrow \text{child}::n, \rho}
\end{array}$$

$$\begin{array}{l}
\text{parent-type}(\rho) = \text{element} * \{ \text{AnyElt}^*, \text{add-anyelt}(\rho), \text{AnyElt}^* \} \\
\text{add-anyelt}(\rho) = () \\
\text{add-anyelt}(\varphi, u) = u \\
\text{add-anyelt}(\rho_1 | \rho_2) = \text{add-anyelt}(\rho_1) | \text{add-anyelt}(\rho_2) \\
\text{add-anyelt}(\rho_1, \rho_2) = \text{add-anyelt}(\rho_1), \text{AnyElt}^*, \text{add-anyelt}(\rho_2) \\
\text{add-anyelt}(\rho^+) = (\text{AnyElt}^*, \text{add-anyelt}(\rho))^+
\end{array}$$

Fig. 14. Inference rules for child.

with the rule FSIBL, the rule FSIBL-FORMULA ensures that if the output type is a single formula type, then the context node has only one following sibling satisfying  $n$ .

#### 4.2.4. Child nodes

As the inference rules for parent are defined in terms of those for self, rules for child can be defined in terms of those for self-fsibl (self or following siblings, a variant of fsibl, defined in the next subsection). As shown in Fig. 13, we first infer a formula  $\varphi$  for self-fsibl and then use it as a constraint for the leftmost child of the context node by adding either  $\langle 1 \rangle$  or  $\langle 1 \rangle$  to  $\varphi$ . Specifically, if the output type is nullable, which means that the context node may not have a child, then we use universal modality (rule CHILD-NULABLE in Fig. 14). Otherwise, the context node always has a child and therefore we use existential modality instead (rule CHILD-NOTNULL in Fig. 14).

In addition, to infer a unit type for the context node, we use an auxiliary function  $\text{parent-type}(\rho)$ , defined in Fig. 14, which computes the type of any node that has some children of type  $\rho$  and possibly more of arbitrary types. To this end, it exploits another auxiliary function  $\text{add-anyelt}(\rho)$  which extracts all unit types at top level of  $\rho$ , while maintaining their order, and adds AnyElt\* between unit types, indicating that there may be more child nodes. Note that  $\text{parent-type}(\rho)$  approximates the type of the context node. For example, consider the tree in Fig. 13. If only nodes A and Y are returned by  $C/\text{child}::n$ , then other nodes such as W, X, and Z must not satisfy the name test  $n$ . This constraint is encoded in the inferred input formula, as discussed in the next subsection, but not in the inferred unit type. If we add negation of a name test, i.e.,  $\neg n$ , we could infer a more precise unit type. However, since all the constraints are already encoded in the inferred input formula, we do not add  $\neg n$  in the definition of regular tree types. Still, the more precise we infer a unit type, more precise we can develop an inference system for XQuery in Section 5, and thus we do not simply use AnyElt in the regular tree type part of the input type.

#### 4.2.5. Self or following siblings

While inference rules for self-fsibl are similar to those for fsibl, there is a key difference. Suppose  $\text{fsibl}::n$  returns nothing (i.e., the output type is  $()$ ). Then it means that there are no following siblings satisfying the name test  $n$ , and thus the input context node should have type  $\mu X.[2](\neg k(n) \wedge X)$  (either there are no following siblings or if any, they do not satisfy  $n$ ). In contrast, if  $\text{self-fsibl}::n$  returns nothing, it means that the context node does not satisfy  $n$  and neither do its following siblings, i.e.,  $\mu X.(\neg k(n) \wedge [2] X)$ .

This difference leads to two interpretations of the output type depending on whether it is nullable or not. To illustrate, assume that the output type is  $(\varphi, u), \rho$ . As in the inference rules for fsibl, we examine the output sequence type from right to left. Suppose that a formula  $\psi$  is inferred from  $\rho$  and that there exists a node  $f$  satisfying  $(\varphi, u)$  ( $f$  can be either the context node or one of its following siblings). Then,  $\psi$  is a constraint on  $f$ 's right next sibling. For example, in Fig. 13,

$$\begin{array}{c}
\text{optional} ::= \text{true} \mid \text{false} \\
\\
\text{SELF-FSIBL} \\
\frac{\varphi \leftarrow \text{self-fsibl}::n, \rho \text{ with } \mu X.(\neg k(n) \wedge [2]X), \text{true}}{\varphi \leftarrow \text{self-fsibl}::n, \rho} \\
\\
\text{SFSIBL-TRUE} \\
\frac{\varphi' = \mu X.(\varphi \wedge k(n) \wedge \text{form}(u) \wedge [2]\psi) \vee (\neg k(n) \wedge \langle 2 \rangle X)}{\varphi' \leftarrow \text{self-fsibl}::n, (\varphi, u) \text{ with } \psi, \text{true}} \\
\\
\text{SFSIBL-FALSE} \\
\frac{\varphi' = \mu X.(\varphi \wedge k(n) \wedge \text{form}(u) \wedge \langle 2 \rangle \psi) \vee (\neg k(n) \wedge \langle 2 \rangle X)}{\varphi' \leftarrow \text{self-fsibl}::n, (\varphi, u) \text{ with } \psi, \text{false}} \\
\\
\text{SFSIBL-SEQ} \\
\frac{\varphi_2 \leftarrow \text{self-fsibl}::n, \rho_2 \text{ with } \psi, \text{optional} \quad \varphi_1 \leftarrow \text{self-fsibl}::n, \rho_1 \text{ with } \varphi_2, \text{optional} \wedge \text{nullable}(\rho_2)}{\varphi_1 \leftarrow \text{self-fsibl}::n, (\rho_1, \rho_2) \text{ with } \psi, \text{optional}}
\end{array}$$

Fig. 15. Inference rules for self-fsibl.

if  $A$  is the context node for `self-fsibl`, then  $f$  should be one of the nodes between  $A$  and  $Z$ . If  $f$  is  $X$ , then  $\psi$  should be true at  $X$ 's right next sibling  $Y$ . Moreover, if  $\rho$  is nullable, then the exact constraint on  $f$  is  $\varphi \wedge [2]\psi$  indicating that  $f$  may not have following siblings. This is the case when  $f$  is  $Z$  in Fig. 13. If  $\rho$  is not nullable, then the exact constraint on  $f$  is  $\varphi \wedge \langle 2 \rangle \psi$  indicating that  $f$  has at least one following sibling and its first following sibling is of type  $\psi$ . Therefore, given an output type  $(\rho_1, \rho_2)$ , when examining  $\rho_1$ , we need to exploit the nullability of  $\rho_2$ .

To this end, we introduce a new judgment  $\varphi \leftarrow \text{self-fsibl}::n, \rho \text{ with } \psi, \text{optional}$  where `optional` denotes either true or false. In this judgment, the meaning of  $\psi$  is twofold: it may denote the constraint on either the context node, i.e., `self`, or one of its following siblings i.e., `fsibl`. The former is when  $\rho$  is  $()$ . For the latter case,  $\psi$  does not refer to the last node in the sequence returned by the axis unlike previous with judgments, but to the one that immediately follows it, which may not exist (thus we use the optional parameter). Suppose that `self-fsibl::n` returns a sequence of nodes  $f_1, \dots, f_n$  of type  $\rho$ . Then  $\psi$  is true at the right next sibling of  $f_n$  where  $f_n$  is the rightmost following sibling returned by `self-fsibl::n` given the output type  $\rho$ . For example, in Fig. 13, if  $A/\text{self-fsibl}::n$  returns  $A, W, X$  and  $\varphi \leftarrow \text{self-fsibl}::n, \rho \text{ with } \psi, \text{optional}$ , then  $\psi$  is true at  $Y$ . This is in clear contrast to the interpretations for other axes where the `with` parameter is true at the lastly reached node among the nodes returned by navigating the given axis. For example, if  $A/\text{fsibl}::n$  returns  $W, X$  and  $\varphi \leftarrow \text{fsibl}::n, \rho \text{ with } \psi$ , then  $\psi$  is true at  $X$ . In the rule SELF-FSIBL in Fig. 15, therefore, the initial `with` parameter is set to  $\mu X.(\neg k(n) \wedge [2]X)$  which means that all the following siblings of the current context node (including the context node itself if the output type is  $()$ ) do not satisfy the name test  $n$ . In addition, the initial `optional` parameter is set to true.

The nullability parameter is examined only when the output type is a single formula type  $(\varphi, u)$ . Consider a judgment  $\varphi' \leftarrow \text{self-fsibl}::n, (\varphi, u) \text{ with } \psi, \text{optional}$ . Then, there should be a node  $f_1$  that satisfies the name test  $n$  and is of type  $(\varphi, u)$  (it can be either the context node or one of its following siblings). Moreover,  $f_1$ 's right next sibling  $f_2$  must satisfy  $\psi$ . If `optional` is true, then  $f_2$  may not exist and thus  $f_1$  has type  $\varphi \wedge k(n) \wedge \text{form}(u) \wedge [2]\psi$  (rule SFSIBL-TRUE). Otherwise,  $f_2$  must exist and thus  $f_1$  has type  $\varphi \wedge k(n) \wedge \text{form}(u) \wedge \langle 2 \rangle \psi$  (rule SFSIBL-FALSE).

For the rest of the cases, we reuse the inference rules in Fig. 7 with minor modifications. For the rules AXIS-EMPTY, AXIS-OR, and AXIS-PLUS, we add one more parameter `optional` in each judgment. Given a sequence type  $(\rho_1, \rho_2)$ , the nullability is updated when examining  $\rho_1$  as shown in the rule SFSIBL-SEQ in Fig. 15. Precisely, the last node  $f$  among the nodes returned by `self-fsibl::n` with the output type  $\rho_1$  may not have following siblings if  $\rho_2$  is nullable and the given `optional` parameter is true. In this case, we use  $[2]\varphi_2$  as a constraint on  $f$  (in combination with the rule SFSIBL-TRUE). Otherwise,  $f$  must have a following sibling and we use  $\langle 2 \rangle \varphi_2$  as a constraint on  $f$  (in combination with the rule SFSIBL-FALSE).

To show that Invariant 4.1 holds for `child`, we need to show that if  $\varphi \leftarrow \text{self-fsibl}::n, \rho$  then either  $[1]\varphi <: \text{parent-type}(\rho)$  or  $\langle 1 \rangle \varphi <: \text{parent-type}(\rho)$  depending on the nullability of  $\rho$ , which is proved by the following lemma.

**Lemma 4.3.** *Suppose  $\varphi \leftarrow \text{self-fsibl}::n, \rho \text{ with } \psi, \text{optional}$ . If `optional` is true, then  $[1]\varphi <: \text{parent-type}(\rho)$ . Otherwise,  $\langle 1 \rangle \varphi <: \text{parent-type}(\rho)$ .*

**Proof.** By induction on a derivation of  $\varphi \leftarrow \text{self-fsibl}::n, \rho \text{ with } \psi, \text{optional}$ . Below we only sketch the proof ideas for the cases where `optional` = false; the proof for the other cases is similar.

Case 1)  $\rho = ()$ : trivial since  $\text{parent-type}(\rho) = \text{AnyElt}$ .

Case 2)  $\rho = \rho_1 \mid \rho_2$ :

(1)  $\varphi_i \leftarrow \text{self-fsibl}::n, \rho_i \text{ with } \psi, \text{false}$  for  $i = 1, 2$

(2)  $\langle 1 \rangle \varphi_i <: \text{parent-type}(\rho_i)$  for  $i = 1, 2$

from the rule AXIS-OR  
by I.H.

(3) Then,  $\langle 1 \rangle (\varphi_1 \vee \varphi_2) <: \text{parent-type}(\rho_1) \mid \text{parent-type}(\rho_2) = \text{parent-type}(\rho_1 \mid \rho_2)$ .

Case 3)  $\rho = (\varphi, u)$ :

- (1)  $\varphi' = \mu X. (\varphi \wedge k(n) \wedge \text{form}(u) \wedge \langle 2 \rangle \psi) \vee (\neg k(n) \wedge \langle 2 \rangle X)$  from the rule SFSIBL-FALSE
- (2)  $\text{parent-type}((\varphi, u)) = \text{element} * \{ \text{AnyElt}^*, u, \text{AnyElt}^* \}$  by definition
- (3)  $\varphi'$  can be thought of as a sequence type  $((\neg k(n))^*, \varphi \wedge k(n) \wedge \text{form}(u) \wedge \langle 2 \rangle \psi)$  using the regular expression notation. Moreover,  $(\neg k(n))^* <: \text{AnyElt}^*$  and  $(\varphi \wedge k(n) \wedge \text{form}(u) \wedge \langle 2 \rangle \psi) <: u, \text{AnyElt}^*$ .
- (4) Therefore,  $\langle 1 \rangle \varphi' <: \text{parent-type}((\varphi, u))$ .

Case 4)  $\rho = \rho_1, \rho_2$ :

- (1)  $\varphi_2 \leftarrow \text{self-fsibl}::n, \rho_2$  with  $\psi$ , false from the rule SFSIBL-SEQ
- (2)  $\varphi_1 \leftarrow \text{self-fsibl}::n, \rho_1$  with  $\varphi_2$ , false from the rule SFSIBL-SEQ
- (3)  $\langle 1 \rangle \varphi_i <: \text{parent-type}(\rho_i)$  for  $i = 1, 2$  by I.H.
- (4) By interpreting  $\varphi_1$  as a type of a sequence of nodes (the context node and its following siblings), we have  $\varphi_1 <: \text{AnyElt}^*, \text{add-anyelt}(\rho_1), \text{AnyElt}^*$ . Similarly,  $\varphi_2 <: \text{AnyElt}^*, \text{add-anyelt}(\rho_2), \text{AnyElt}^*$ .
- (5) Moreover, since  $\varphi_1$  contains  $\varphi_2$  as the rightmost subformula (the one that the most  $\langle 2 \rangle$ s precede directly or indirectly through recursion), we have  $\varphi_1 <: \text{AnyElt}^*, \text{add-anyelt}(\rho_1), \text{AnyElt}^*, \text{add-anyelt}(\rho_2), \text{AnyElt}^*$ .
- (6) Therefore,  $\langle 1 \rangle \varphi_1 <: \text{parent-type}(\rho_1, \rho_2)$ .

Case 5)  $\rho = \rho_1^+$ :

- (1)  $\varphi \leftarrow \text{self-fsibl}::n, \rho_1$  with  $X \vee \psi$ , false from the rule AXIS-PLUS
- (2)  $\langle 1 \rangle \varphi <: \text{parent-type}(\rho_1)$  by I.H.
- (3) By interpreting  $\varphi$  as a type of a sequence of nodes, we have  $\varphi <: \text{AnyElt}^*, \text{add-anyelt}(\rho_1), \text{AnyElt}^*$ , where the rightmost subformula of  $\varphi$ , i.e.,  $X \vee \psi$ , is subsumed by the second  $\text{AnyElt}^*$ .
- (4) Therefore,  $\mu X. \varphi <: (\text{AnyElt}^*, \text{add-anyelt}(\rho_1))^+, \text{AnyElt}^*$  and thus  $\langle 1 \rangle \mu X. \varphi <: \text{parent-type}(\rho_1^+)$ .  $\square$

#### 4.2.6. Descendants

`desc` is more complicated than other axes because its semantics does not have a “closure property” with respect to the document order. Formally, we say that *axis* is “forward-closed” if  $\llbracket f/\text{axis}::n \rrbracket = f_1, \dots, f_n$  implies  $\llbracket f_i/\text{axis}::n \rrbracket = f_{i+1}, \dots, f_n$  for any  $i$ . Similarly, *axis* is “backward-closed” if  $\llbracket f/\text{axis}::n \rrbracket = f_1, \dots, f_n$  implies  $\llbracket f_i/\text{axis}::n \rrbracket = f_1, \dots, f_{i-1}$  for any  $i$ . Note that `fsibl` is forward-closed while `psibl` and `anc` are backward-closed. Although `desc` itself is not forward-closed, it can be defined in terms of another forward-closed, parameterized axis `desc-or-foll( $f$ )` which returns the descendants and following nodes of the context node that appear in document order before the first following node of  $f$ . We remark here that in the XML terminology the following nodes of  $f$  mean the nodes that appear after  $f$  in document order but are not a descendant of  $f$ . More precisely, if  $\llbracket f/\text{desc}::n \rrbracket = f_1, \dots, f_n$ , then  $\llbracket f_i/\text{desc-or-foll}(f)::n \rrbracket = f_{i+1}, \dots, f_n$  for any  $i$ . Note that  $f_{i+1}$  is  $f_i$ 's descendant or one of its following nodes. Below, based on this observation, we develop inference rules for `desc` using a judgment of the form  $\varphi \leftarrow \text{desc}::n, \rho$  with  $\psi$  where the with parameter  $\psi$  now denotes a constraint on the last node in document order among the descendants and following nodes of the context node returned by `desc-or-foll` when the output type is  $\rho$ .

In the rule `DESC` in Fig. 16, the initial with parameter is much more complicated than other axes because we need to specify constraints only on the descendants of the context node, but not on others. More precisely, if the output type is a sequence type of the form  $(\varphi_1, u_1), \dots, (\varphi_n, u_n)$ , then  $f/\text{desc}::n$  for any input node  $f$  returns a sequence  $f_1, \dots, f_n$  of descendants in document order, each of which has type  $(\varphi_i, u_i)$ . Now, to infer an exact input type, we need to specify that all the nodes that follow  $f_n$  but precede  $f$ 's first following node must not satisfy the name test  $n$ . This constraint is encoded in `noNextUpTo( $k(n), \alpha$ )` which exploits a nominal denoted by  $\alpha$ . A nominal is simply an atomic proposition like node labels  $\sigma$  but it holds at exactly one node of any given tree [28]. Specifically, in the rule `DESC`, the nominal  $\alpha$  is true only at the context node on which `desc` is applied, which is guaranteed by `noWhereElse( $\alpha$ )` in the final input type. This nominal  $\alpha$  is used as a search bound for descendants during the inference process, i.e., `noNextUpTo( $k(n), \alpha$ )`.

If the output type is a single formula type  $(\varphi, u)$  and the with parameter is  $\psi$ , it means that the context node has a descendant  $f$  of type  $(\varphi, u)$  that satisfies the name test  $n$  and at which  $\psi$  is true (rule `DESC-FORMULA`). Moreover, any node between the context node and  $f$  in document order should not satisfy the name test  $n$ . All these constraints are encoded in `fstDescFoll( $\varphi \wedge k(n) \wedge \text{form}(u) \wedge \psi, k(n)$ )`, defined in Fig. 16, which specifies the first node satisfying  $n$  among the descendants and following nodes of the context node. As for other cases, we simply use the inference rules in Fig. 7. In particular, the rule `AXIS-FORWARD-SEQ` can be used as it is since each descendant is connected to the next descendant in document order by `fstDescFoll()` (rule `DESC-FORMULA`) and the last descendant in document order satisfies `noNextUpTo()` for the input node on which `desc` is initially applied (rule `DESC`).



$$\text{DESC} \frac{\varphi \leftarrow \text{desc} : : n, \rho \text{ with } \text{noNextUpTo}(k(n), \alpha)}{(\varphi \wedge \text{noWhereElse}(\alpha), \text{AnyElt}) \leftarrow \text{desc} : : n, \rho} \quad (\alpha \text{ fresh})$$

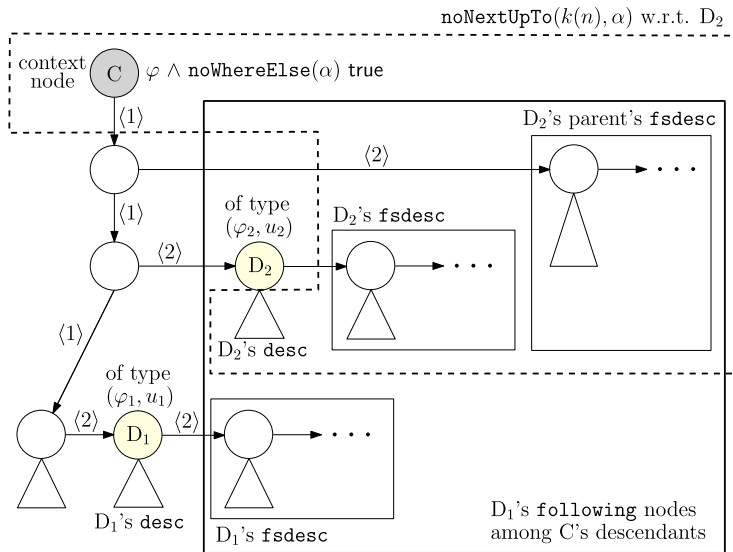
$$\text{DESC-FORMULA} \frac{\varphi' = \text{fstDescFoll}(\varphi \wedge k(n) \wedge \text{form}(u) \wedge \psi, k(n))}{\varphi' \leftarrow \text{desc} : : n, (\varphi, u) \text{ with } \psi}$$

**Auxiliary definitions:**

$$\begin{aligned} \chi ? \psi_1 : \psi_2 &\equiv (\chi \wedge \psi_1) \vee (\neg \chi \wedge \psi_2) \\ \text{has-desc}(\chi) &= \langle 1 \rangle (\mu Z. \chi \vee \langle 1 \rangle Z \vee \langle 2 \rangle Z) \\ \text{has-fsdesc}(\chi) &= \langle 2 \rangle (\mu Z. \chi \vee \langle 1 \rangle Z \vee \langle 2 \rangle Z) \\ \text{has-prec}(\chi) &= \mu Z. \langle \bar{1} \rangle Z \vee \langle \bar{2} \rangle (\chi \vee \text{has-desc}(\chi) \vee Z) \\ \text{has-foll}(\chi) &= \mu Z. \text{has-fsdesc}(\chi) \vee \text{has-parent}(Z) \\ \text{noNextUpTo}(\chi, \alpha) &= \neg \text{has-desc}(\chi) \wedge \mu Z. \alpha ? \top : (\neg \text{has-fsdesc}(\chi) \wedge \text{has-parent}(Z)) \\ \text{noWhereElse}(\chi) &= \chi \wedge \neg (\text{has-anc}(\chi) \vee \text{has-prec}(\chi) \vee \text{has-desc}(\chi) \vee \text{has-foll}(\chi)) \\ \text{fstSelfFsDesc}(\chi_1, \chi_2) &= \mu Z. \chi_1 \vee (\neg \chi_2 \wedge (\text{has-desc}(\chi_2) ? \langle 1 \rangle Z : \langle 2 \rangle Z)) \\ \text{fstFoll}(\chi_1, \chi_2) &= \mu Z. \langle 2 \rangle \text{fstSelfFsDesc}(\chi_1, \chi_2) \vee (\neg \text{has-fsdesc}(\chi_2) \wedge \text{has-parent}(Z)) \\ \text{fstDescFoll}(\chi_1, \chi_2) &= \langle 1 \rangle \text{fstSelfFsDesc}(\chi_1, \chi_2) \vee (\neg \text{has-desc}(\chi_2) \wedge \text{fstFoll}(\chi_1, \chi_2)) \end{aligned}$$

- $\text{has-desc}(\chi)$ : there is a descendant satisfying  $\chi$ .
- $\text{has-fsdesc}(\chi)$ : there is a node satisfying  $\chi$  among the following siblings and their descendants.
- $\text{has-prec}(\chi)$ : there is a node satisfying  $\chi$  which precedes the context node in document order and is not an ancestor.
- $\text{has-foll}(\chi)$ : there is a node satisfying  $\chi$  which follows the context node in document order and is not a descendant.
- $\text{noNextUpTo}(\chi, \alpha)$ : there is no node satisfying  $\chi$  which appears strictly after the context node in document order and before a node satisfying  $\alpha$  (invariant:  $\alpha$  should denote a nominal).
- $\text{noWhereElse}(\chi)$ : only the context node satisfies  $\chi$ .
- $\text{fstSelfFsDesc}(\chi_1, \chi_2)$ : the first node  $f$  in document order of the set {self, all following siblings, and all their descendants} that satisfies  $\chi_1$ . Any node preceding  $f$  in the set does not satisfy  $\chi_2$ .
- $\text{fstFoll}(\chi_1, \chi_2)$ : the first node  $f$  satisfying  $\chi_1$  among the nodes reachable by navigating following. Any node between the context node and  $f$  reached by navigating following does not satisfy  $\chi_2$ .
- $\text{fstDescFoll}(\chi_1, \chi_2)$ : the first node  $f$  satisfying  $\chi_1$  that appears strictly after the context node in document order. Any node between the context node and  $f$  does not satisfy  $\chi_2$ .
- Note. In  $\text{fstSelfFsDesc}(\chi_1, \chi_2)$ ,  $\text{fstFoll}(\chi_1, \chi_2)$ , and  $\text{fstDescFoll}(\chi_1, \chi_2)$ , we assume  $\langle \chi_1 \rangle \subseteq \langle \chi_2 \rangle$ . This is always ensured by the rule DESC-FORMULA which is the only rule that uses these auxiliary functions.

**Fig. 16.** Inference rules for desc.



**Fig. 17.**  $\varphi \leftarrow \text{desc} : : n, ((\varphi_1, u_1), (\varphi_2, u_2))$  with  $\text{noNextUpTo}(k(n), \alpha)$  where C is the context node and D<sub>1</sub> and D<sub>2</sub> are the only nodes satisfying the name test  $n$ , each of which has type  $(\varphi_1, u_1)$  and  $(\varphi_2, u_2)$ , respectively. D<sub>1</sub> precedes D<sub>2</sub> in document order and D<sub>2</sub> is the first node satisfying  $n$  among D<sub>1</sub>'s descendants and following nodes.

To illustrate, consider an example tree in Fig. 17. Suppose that the output type is  $(\varphi_1, u_1), (\varphi_2, u_2)$ . Then there exist only two descendants satisfying the name test  $n$ , namely,  $D_1$  and  $D_2$ . According to the rule **AXIS-FORWARD-SEQ**, we first analyze the rightmost output type  $(\varphi_2, u_2)$ . In other words, we first infer a constraint on the node  $D_2$ . Since  $D_2$  is the last node returned by `desc : n, noNextUpTo(k(n),  $\alpha$ )` should be true at  $D_2$  which means that  $k(n)$  is not true at  $D_2$ 's descendants, its following siblings and their descendants, its parent's following siblings and their descendants, its parent's parent's following siblings and their descendants, and so on until the initial context node  $C$ , marked with a nominal  $\alpha$ , is reached (rule **DESC**). Moreover, when locally analyzing  $D_2$  with the output type  $(\varphi_2, u_2)$ , the context node is  $D_1$ . From  $D_1$ 's perspective,  $D_2$  is the first node satisfying the name test  $n$  among  $D_1$ 's descendants and following nodes. This constraint is expressed by using the function `fstDescFollow()` (rule **DESC-FORMULA**).

### 4.3. Properties of backward type inference for XPath axes

In this section, we briefly discuss the soundness and completeness of our backward type inference system for XPath axes. In other words, our backward inference is exact.

**Theorem 4.4** (Exact type inference). *Suppose  $\rho_i \leftarrow \text{axis} : : n, \rho_o$ . Then,  $f \in \llbracket \rho_i \rrbracket$  if and only if  $\llbracket f / \text{axis} : : n \rrbracket \in \llbracket \rho_o \rrbracket$ .*

In Theorem 4.4, the only-if-direction states the soundness and the if-direction states the completeness. More precisely, the soundness states that given an axis  $\text{axis} : : n$  and an output type  $\rho_o$ , if our inference system infers an input type  $\rho_i$  and some focused tree is of type  $\rho_i$ , then it always produces a sequence of nodes of type  $\rho_o$ . In contrast, the completeness states the opposite, that is, given an axis  $\text{axis} : : n$ , if some focused tree  $f$  produces a sequence of nodes of type  $\rho_o$ , then our inference system can infer an input type  $\rho_i$  from  $\rho_o$  and  $f$  is of type  $\rho_i$ . In particular, our inference system fails only if the axis is `self : : n` or `parent : : n` and the output type is a sequence type  $\rho_1, \rho_2$  where both  $\rho_1$  and  $\rho_2$  are not nullable. Note that this case never happens though since there is only one self and parent, if any. To prove Theorem 4.4, we use Proposition 4.2 which is proved by Lemmas 4.5 and 4.6 for the auxiliary judgment  $\varphi \leftarrow \text{axis} : : n, \rho$  with  $\psi$ .

**Lemma 4.5** (Soundness). *Suppose  $\varphi \leftarrow \text{axis} : : n, \rho$  with  $\psi, f \in \langle\langle \varphi \rangle\rangle$ , and  $\llbracket f / \text{axis} : : n \rrbracket = f_1, \dots, f_n$ .*

- For backward axes, let  $f = f_{n+1}$ . Then,  $\exists 1 \leq i \leq n + 1$  such that  $f_i \in \langle\langle \psi \rangle\rangle$  and  $f_i, \dots, f_n \in \llbracket \rho \rrbracket$ .
- For forward axes, let  $f = f_0$ . Then,  $\exists 0 \leq i \leq n$  such that  $f_i \in \langle\langle \psi \rangle\rangle$  and  $f_1, \dots, f_i \in \llbracket \rho \rrbracket$ .

Lemma 4.5 is a one-way formalization of the interpretation of the judgment  $\varphi \leftarrow \text{axis} : : n, \rho$  with  $\psi$ . To illustrate, consider Fig. 5 again. In the figure,  $f_{n+1} = C$  and  $f_i = A$  where  $C \in \langle\langle \varphi \rangle\rangle$  and  $A \in \langle\langle \psi \rangle\rangle$ . Moreover, the sequence  $f_i, \dots, f_n$  of nodes selected from the node  $A$  to the node  $B$  by `psibl : : n` has type  $\rho$ . Below we show some cases of the proof of Lemma 4.5.

**Proof of Lemma 4.5.** By induction on a derivation of  $\varphi \leftarrow \text{axis} : : n, \rho$  with  $\psi$ . Here, we show some cases where  $\text{axis}$  is a backward axis.

Case 1)  $\rho = ()$ :

- (1)  $\varphi = \psi$  from the rule **AXIS-EMPTY**
- (2) Let  $i = n + 1$ .
- (3) Then,  $f_i \in \langle\langle \psi \rangle\rangle$  and  $f_i, \dots, f_n = \epsilon \in \llbracket () \rrbracket$ . from assumptions

Case 2)  $\rho = \rho_1, \rho_2$ :

- (1)  $\varphi_1 \leftarrow \text{axis} : : n, \rho_1$  with  $\psi$  from the rule **AXIS-BACKWARD-AXIS**
- (2)  $\varphi \leftarrow \text{axis} : : n, \rho_2$  with  $\varphi_1$  from the rule **AXIS-BACKWARD-AXIS**
- (3)  $\exists 1 \leq j \leq n + 1$  s.t.  $f_j \in \langle\langle \varphi_1 \rangle\rangle$  and  $f_j, \dots, f_n \in \llbracket \rho_2 \rrbracket$  by I.H. on (2)
- (4)  $\llbracket f_j / \text{axis} : : n \rrbracket = f_1, \dots, f_{j-1}$
- (5)  $\exists 1 \leq k \leq j$  s.t.  $f_k \in \langle\langle \psi \rangle\rangle$  and  $f_k, \dots, f_{j-1} \in \llbracket \rho_1 \rrbracket$  by I.H. on (1) and (4)
- (6) Let  $i = k$ .
- (7) Then,  $f_i \in \langle\langle \psi \rangle\rangle$  and  $f_i, \dots, f_{j-1}, f_j, \dots, f_n \in \llbracket (\rho_1, \rho_2) \rrbracket$ . □

Below we state the completeness lemma for the auxiliary judgment. In particular, given an axis  $\text{axis} : : n$  and an output type  $\rho$ , we assume that our inference system always infers some input formula  $\varphi$  since the inference never fails. (The inferred formula may be unsatisfiable though, indicating a contradiction.)

**Lemma 4.6** (Completeness). *Suppose  $\llbracket f / \text{axis} : : n \rrbracket = f_1, \dots, f_n$  and  $\varphi \leftarrow \text{axis} : : n, \rho$  with  $\psi$ .*

- For backward axes, let  $f = f_{n+1}$ . If  $\exists 1 \leq i \leq n + 1$  such that  $f_i \in \langle\langle \psi \rangle\rangle$  and  $f_1, \dots, f_n \in \llbracket \rho \rrbracket$ , then  $f \in \langle\langle \varphi \rangle\rangle$ .
- For forward axes, let  $f = f_0$ . If  $\exists 0 \leq i \leq n$  such that  $f_i \in \langle\langle \psi \rangle\rangle$  and  $f_1, \dots, f_n \in \llbracket \rho \rrbracket$ , then  $f \in \langle\langle \varphi \rangle\rangle$ .

**Proof.** By induction on a derivation of  $\varphi \leftarrow \text{axis} : : n, \rho$  with  $\psi$ . Here, we only show the case where  $\rho = \rho_1, \rho_2$  and axis is a backward axis.

(1) $\varphi_1 \leftarrow \text{axis} : : n, \rho_1$ with $\psi$	from the rule AXIS-BACKWARD-AXIS
(2) $\varphi \leftarrow \text{axis} : : n, \rho_2$ with $\varphi_1$	from the rule AXIS-BACKWARD-AXIS
(3) $\exists i \leq j \leq n$ s.t. $f_i, \dots, f_{j-1} \in \llbracket \rho_1 \rrbracket$ and $f_j, \dots, f_n \in \llbracket \rho_2 \rrbracket$	from $f_i, \dots, f_n \in \llbracket (\rho_1, \rho_2) \rrbracket$
(4) $\llbracket f_j / \text{axis} : : n \rrbracket = f_1, \dots, f_{j-1}$	
(5) $1 \leq i \leq j$ and $f_i \in \langle\langle \psi \rangle\rangle$	from (3) and assumptions
(6) $f_j \in \langle\langle \varphi_1 \rangle\rangle$	by I.H. on (1), (3), (4), (5)
(7) $f_{n+1} \in \langle\langle \varphi \rangle\rangle$	by I.H. on (2), (3), (6) $\square$

## 5. Inference for the XQuery core

In this section, we present our backward type inference system for the XQuery core in the style of constraint solving systems [29–31], building on the results of the previous section. We first clarify what we infer from the given expression  $e$  and output type  $\rho$ . Precisely, we use a judgment of the form  $\mathcal{S} \leftarrow e : \rho$  which means that given an expression  $e$  and an output type  $\rho$ , it generates a set  $\mathcal{S}$  of constraint-sets for free variables in  $e$  where free and bound variables are defined in the usual way. Our goal is then to design inference rules that ensure that if we substitute those free variables with any sequences of focused trees satisfying one of the constraint-sets in  $\mathcal{S}$ ,  $e$  evaluates to a value, *i.e.*, a sequence of focused trees, that has the type  $\rho$ . By convention, if  $\mathcal{S}$  is an empty set, it is unsatisfiable, and we denote it by  $\emptyset$ . In contrast, a singleton set consisting of an empty set is always satisfiable, and we denote it by  $\perp$ .

Formally, a constraint-set  $C$  is a set of bindings of variables with formula-enriched sequence types, where each binding is denoted by  $(\$var : \rho)$ . Given a constraint-set  $C$ , we consider any for-loop and let-bound variables not appearing in  $C$  to be implicitly bound to  $(\top, \text{AnyElt})$  and  $(\top, \text{AnyElt})^*$ , respectively. Moreover, a constraint-set  $C$  is unsolvable if it contains a constraint specifying that a variable should satisfy  $\perp$ , for example,  $(\$var : (\perp, u))$  or  $(\$var : (\varphi, u))$  where  $\varphi$  is unsatisfiable. We simply write  $\{\perp\}$  to denote such an unsolvable constraint-set. If  $\mathcal{S}$  contains  $\{\perp\}$ , we can safely remove it from  $\mathcal{S}$ . We often consider a constraint-set  $C$  to be a mapping from variables to their types and thus use the usual notations such as:

$$\begin{aligned}
 \text{dom}(C) &\stackrel{\text{def}}{=} \{\$var \mid (\$var : \rho) \in C\} \\
 C(\$var) &\stackrel{\text{def}}{=} \rho && \text{if } (\$var : \rho) \in C \\
 C(\$v) &\stackrel{\text{def}}{=} (\top, \text{AnyElt}) && \text{if } \$v \notin \text{dom}(C) \\
 C(\$v) &\stackrel{\text{def}}{=} (\top, \text{AnyElt})^* && \text{if } \$v \notin \text{dom}(C)
 \end{aligned}$$

We also introduce the following operations.

**Definition 5.1.** Let  $C_1$  and  $C_2$  be constraint-sets, which are not  $\{\perp\}$ , and  $\mathcal{S}, \mathcal{S}_1$ , and  $\mathcal{S}_2$  be sets of constraint-sets. We define:

$$\begin{aligned}
 C_1 \sqcap C_2 &\stackrel{\text{def}}{=} \{(\$var : \rho) \in C_1 \mid \$var \notin \text{dom}(C_2)\} \cup \\
 &\quad \{(\$var : \rho) \in C_2 \mid \$var \notin \text{dom}(C_1)\} \cup \\
 &\quad \{(\$var : \rho_1 \wedge \rho_2) \mid (\$var : \rho_1) \in C_1 \text{ and } (\$var : \rho_2) \in C_2\} \\
 C \setminus \$var_0 &\stackrel{\text{def}}{=} \{(\$var : \rho) \in C \mid \$var \neq \$var_0\} \\
 \mathcal{S}_1 \sqcap \mathcal{S}_2 &\stackrel{\text{def}}{=} \{C_1 \sqcap C_2 \mid C_1 \in \mathcal{S}_1, C_2 \in \mathcal{S}_2\} \\
 \mathcal{S}_1 \sqcup \mathcal{S}_2 &\stackrel{\text{def}}{=} \mathcal{S}_1 \cup \mathcal{S}_2 \\
 \mathcal{S} \setminus \$var &\stackrel{\text{def}}{=} \{C \setminus \$var \mid C \in \mathcal{S}\}
 \end{aligned}$$

For any constraint-set  $C$ ,  $C \sqcap \{\perp\} = \{\perp\} \sqcap C = \{\perp\}$ .

In the definition above, we use  $\rho_1 \wedge \rho_2$  to denote the intersection of  $\rho_1$  and  $\rho_2$  whose semantics  $\llbracket \rho_1 \wedge \rho_2 \rrbracket$  is inductively defined as  $\llbracket \rho_1 \rrbracket \cap \llbracket \rho_2 \rrbracket$ . In other words, for any focused tree  $f$ ,  $f \in \llbracket \rho_1 \wedge \rho_2 \rrbracket$  if and only if  $f \in \llbracket \rho_1 \rrbracket$  and  $f \in \llbracket \rho_2 \rrbracket$ . Although we use intersection types only internally during type inference, they can seamlessly be added into the external language [32].

$$\begin{array}{c}
\text{I-OR} \\
\frac{\mathcal{S}_i \leftarrow e : \rho_i}{\mathcal{S}_1 \sqcup \mathcal{S}_2 \leftarrow e : \rho_1 \mid \rho_2} \quad (i = 1, 2) \\
\\
\text{I-FVAR} \\
\frac{\rho' \leftarrow \text{self} : *, \rho}{\{(\$v : \rho')\} \leftarrow \$v : \rho} \\
\\
\text{I-LVAR} \\
\frac{}{\{(\$v : \rho)\} \leftarrow \$v : \rho} \\
\\
\text{I-AXIS} \\
\frac{\rho' \leftarrow \text{axis} : n, \rho}{\{(\$v : \rho')\} \leftarrow \$v/\text{axis} : n : \rho} \\
\\
\text{I-ELEMENT} \\
\frac{(\varphi_{\text{root}}, u) <: \rho \quad u = \text{element } n \{ \tau \} \quad \sigma = n \text{ or } \sigma = * \quad \mathcal{S} \leftarrow e : \text{form-enriched}(\tau)}{\mathcal{S} \leftarrow \langle \sigma \rangle \{ e \} \langle / \sigma \rangle : u : \rho} \\
\\
\text{I-IFNONEMPTY} \\
\frac{\mathcal{S}_1 \leftarrow e_1 : (\top, \text{AnyElt})^+ \quad \mathcal{S}_2 \leftarrow e_2 : \rho}{\mathcal{S}_1 \sqcap \mathcal{S}_2 \leftarrow \text{if empty}(e_1) \text{ then } e_2 \text{ else } e_3 : \rho} \\
\\
\text{I-IFEMPTY} \\
\frac{\mathcal{S}_1 \leftarrow e_1 : () \quad \mathcal{S}_3 \leftarrow e_3 : \rho}{\mathcal{S}_1 \sqcap \mathcal{S}_3 \leftarrow \text{if empty}(e_1) \text{ then } e_2 \text{ else } e_3 : \rho} \\
\\
\text{I-IFANY} \\
\frac{\mathcal{S}_1 \leftarrow e_1 : (\top, \text{AnyElt})^* \quad \mathcal{S}_2 \leftarrow e_2 : \rho \quad \mathcal{S}_3 \leftarrow e_3 : \rho}{\mathcal{S}_1 \sqcap \mathcal{S}_2 \sqcap \mathcal{S}_3 \leftarrow \text{if empty}(e_1) \text{ then } e_2 \text{ else } e_3 : \rho} \\
\\
\text{I-LET} \\
\frac{\mathcal{S}_2 \leftarrow e_2 : \rho \quad S = \{ \mathcal{S}_1 \sqcap \{ C \}_{\$v} \mid \mathcal{S}_1 \leftarrow e_1 : C(\$v), C \in \mathcal{S}_2 \}}{\bigsqcup_{\mathcal{S} \in S} \mathcal{S} \leftarrow \text{let } \$v := e_1 \text{ return } e_2 : \rho}
\end{array}$$

#### Auxiliary definitions:

$$\begin{array}{l}
\text{form-enriched}(\ ()) = \ () \\
\text{form-enriched}(u) = \ (\text{form}(u), u) \\
\text{form-enriched}(\tau_1, \tau_2) = \ \text{form-enriched}(\tau_1), \text{form-enriched}(\tau_2) \\
\text{form-enriched}(\tau_1 \mid \tau_2) = \ \text{form-enriched}(\tau_1) \mid \text{form-enriched}(\tau_2) \\
\text{form-enriched}(\tau^+) = \ \text{form-enriched}(\tau)^+ \\
\\
\text{single}(\ ()) = \ \perp \\
\text{single}(\varphi, u) = \ \varphi \wedge \text{form}(u) \\
\text{single}(\rho_1, \rho_2) = \ \begin{cases} \perp & \text{if } \neg \text{nullable}(\rho_1) \text{ and } \neg \text{nullable}(\rho_2) \\ \text{single}(\rho_1) & \text{if } \neg \text{nullable}(\rho_1) \text{ and } \text{nullable}(\rho_2) \\ \text{single}(\rho_2) & \text{if } \text{nullable}(\rho_1) \text{ and } \neg \text{nullable}(\rho_2) \\ \text{single}(\rho_1) \vee \text{single}(\rho_2) & \text{if } \text{nullable}(\rho_1) \text{ and } \text{nullable}(\rho_2) \end{cases} \\
\text{single}(\rho_1 \mid \rho_2) = \ \text{single}(\rho_1) \vee \text{single}(\rho_2) \\
\text{single}(\rho_1 \wedge \rho_2) = \ \text{single}(\rho_1) \wedge \text{single}(\rho_2) \\
\text{single}(\rho^+) = \ \text{single}(\rho)
\end{array}$$

Fig. 18. Backward type inference rules for the XQuery core.

### 5.1. Inference rules

Figs. 18 and 19 show our backward type inference rules for the XQuery core. We first describe the case where the output type is a union type  $\rho_1 \mid \rho_2$  (rule I-OR). In this case, the input constraint is a union of  $\mathcal{S}_1$  and  $\mathcal{S}_2$ , which are inferred from  $\rho_1$  and  $\rho_2$ , respectively. If one of  $\mathcal{S}_i$  is unsatisfiable, i.e.,  $\emptyset$ , it is simply ignored since  $\mathcal{S} \sqcup \emptyset = \mathcal{S}$  for any  $\mathcal{S}$ . If both  $\mathcal{S}_1$  and  $\mathcal{S}_2$  are unsatisfiable, the input constraint is  $\emptyset$  which means that expression  $e$  can never have the output type  $\rho_1 \mid \rho_2$  in the first place. Similarly, if the output type is an intersection type  $\rho_1 \wedge \rho_2$ , the input constraint is an intersection of  $\mathcal{S}_1$  and  $\mathcal{S}_2$ , each of which is inferred from  $\rho_i$  (rule I-AND). In this case, if one of  $\mathcal{S}_i$  is unsatisfiable, then the input type is also unsatisfiable. During the inference, either the rule I-OR or the rule I-AND should first be tried.

The I-EMP, I-FVAR, I-LVAR, and I-AXIS rules are relatively easy. First, in the rule I-EMP, if the output type  $\rho$  is nullable, then the input constraint is  $\mathbb{1}$  which means that  $\epsilon$  is of type  $\rho$  without further constraints. In the rule I-FVAR, we use the inference rules for the `self` axis since a for-loop variable is bound only to an XML element, not a sequence. In contrast, the rule I-LVAR just binds a let-bound variable to the given sequence type since it can be bound to an arbitrary sequence. The rule I-AXIS uses the inference rules for the axis expression, and binds the for-loop variable to the inferred type.

In the rule I-ELEMENT, we consider only a type-annotated element constructor  $\langle \sigma \rangle \{ e \} \langle / \sigma \rangle : u$ . The annotated type  $u$  should be a subtype of the output type  $\rho$  since we are using a backward type inference. Specifically, since an element constructor always reduces to a root element, we check the subtype relation  $(\varphi_{\text{root}}, u) <: \rho$  where  $\varphi_{\text{root}} = \neg(\bar{1}) \top \wedge \neg(\bar{2}) \top \wedge \neg(\bar{2}) \top$  specifies that the given node is a root (the subtype relation is explained shortly). Let  $u$  be `element`  $n \{ \tau \}$ . Then, the node label  $\sigma$  should match the name test  $n$ . Finally, we infer input constraints from the body expression  $e$ , which reduces to a sequence of child nodes, and the output type `form-enriched`( $\tau$ ), which is the type of the child nodes. `form-enriched`( $\tau$ ) enriches the given regular tree type  $\tau$  by simply associating each unit type  $u$  that appears in  $\tau$  with an equivalent downward-only formula `form`( $u$ ), i.e., without context information.

$$\begin{array}{c}
\text{I-SEQ} \\
\frac{S = \{\mathcal{S}_1 \sqcap \mathcal{S}_2 \mid \mathcal{S}_i \leftarrow e_i : \rho_i, (\rho_1, \rho_2) \in \text{split}(\rho)\}}{\bigsqcup_{\mathcal{S} \in S} \mathcal{S} \leftarrow (e_1, e_2) : \rho} \quad \text{I-ERR} \\
\frac{\text{(if no other rule applies)}}{\emptyset \leftarrow e : \rho} \\
\\
\text{I-FOREMPTY} \\
\frac{\mathcal{S}_2 \leftarrow e_2 : () \quad S = \{\mathcal{S}_1 \sqcap \{C \setminus \$v\} \mid \mathcal{S}_1 \leftarrow e_1 : C(\$v)^*, C \in \mathcal{S}_2\}}{\bigsqcup_{\mathcal{S} \in S} \mathcal{S} \sqcup \mathcal{S}_0 \leftarrow \text{for } \$v \text{ in } e_1 \text{ return } e_2 : ()} \quad \mathcal{S}_0 \leftarrow e_1 : () \\
\\
\text{I-FOR} \\
\frac{\mathcal{S} \leftarrow e_2 : \rho \quad \mathcal{S}' \leftarrow e_2 : () \quad \neg \text{nullable}(\rho)}{S = \{\mathcal{S}'' \sqcap \{C \setminus \$v\} \mid \mathcal{S}'' \leftarrow e_1 : C(\$v). \text{Quant}(\rho), C \in \mathcal{S}\}} \\
\mathcal{S}' = \{\mathcal{S}'' \sqcap \{C \setminus \$v \sqcap C' \setminus \$v\} \mid \mathcal{S}'' \leftarrow e_1 : (C'(\$v)^*, C(\$v), C'(\$v)^*). \text{Quant}(\rho), (C, C') \in \mathcal{S} \times \mathcal{S}'\} \quad (\rho \neq ()) \\
\frac{\bigsqcup_{\mathcal{S} \in S \cup \mathcal{S}'} \mathcal{S} \leftarrow \text{for } \$v \text{ in } e_1 \text{ return } e_2 : \rho}{} \\
\\
\text{I-FORNULLABLE} \\
\frac{\mathcal{S} \leftarrow e_2 : \rho \quad \mathcal{S}' \leftarrow e_2 : () \quad \text{nullable}(\rho) \quad \mathcal{S}_0 \leftarrow e_1 : ()}{S = \{\mathcal{S}'' \sqcap \{C \setminus \$v\} \mid \mathcal{S}'' \leftarrow e_1 : C(\$v). \text{Quant}(\rho), C \in \mathcal{S}\}} \\
\mathcal{S}' = \{\mathcal{S}'' \sqcap \{C \setminus \$v \sqcap C' \setminus \$v\} \mid \mathcal{S}'' \leftarrow e_1 : (C'(\$v)^*, C(\$v), C'(\$v)^*). \text{Quant}(\rho), (C, C') \in \mathcal{S} \times \mathcal{S}'\} \quad (\rho \neq ()) \\
\frac{\bigsqcup_{\mathcal{S} \in S \cup \mathcal{S}'} \mathcal{S} \sqcup \mathcal{S}_0 \leftarrow \text{for } \$v \text{ in } e_1 \text{ return } e_2 : \rho}{}
\end{array}$$

**Auxiliary definitions:**

$$\begin{aligned}
\text{split}(\rho) &= \{(\rho_1, \rho_2)\} \\
\text{split}((\varphi, u)) &= \{(\rho_1, (\varphi, u)), ((\varphi, u), \rho_2)\} \\
\text{split}(\rho_1 \mid \rho_2) &= \text{split}(\rho_1) \cup \text{split}(\rho_2) \\
\text{split}(\rho^+) &= \{(\rho_1, \rho_2^+), (\rho_1^+, \rho_2), (\rho_1^+, \rho_2^+)\} \cup \{(\rho_1^*, \rho_1), (\rho_2, \rho_2^*) \mid (\rho_1, \rho_2) \in \text{split}(\rho)\} \\
\text{split}(\rho_1, \rho_2) &= \{(\rho_1, \rho_2)\} \cup \{(\rho_{11}, \rho_{12}, \rho_{22}) \mid (\rho_{11}, \rho_{12}) \in \text{split}(\rho_1)\} \cup \\
&\quad \{(\rho_{21}, \rho_{21}), \rho_{22}\} \mid (\rho_{21}, \rho_{22}) \in \text{split}(\rho_2)\} \\
\text{Quant}(\rho) &= + \quad \text{if } \rho \text{ is of the form } \rho'^+ \quad \rho \cdot + = \rho'^+ \\
\text{Quant}(\rho) &= 1 \quad \text{otherwise} \quad \rho \cdot 1 = \rho
\end{aligned}$$

**Fig. 19.** Backward type inference rules for the XQuery core, continued.

To check the subtype relation  $(\varphi_{\text{root}}, u) <: \rho$ , we first compute the type  $\rho'$  for the set of all single focused tree nodes that are contained in  $\llbracket \rho \rrbracket$ . Then  $(\varphi_{\text{root}}, u) <: \rho$  if  $(\varphi_{\text{root}}, u) <: \rho'$  because  $(\varphi_{\text{root}}, u)$  denotes a set of focused tree nodes. Next, we translate  $u$  and  $\rho'$  into equivalent formulas  $\varphi$  and  $\psi$ , respectively, and then test the satisfiability of  $\varphi_{\text{root}} \wedge \varphi \wedge \neg \psi$ . To this end, we use an auxiliary function  $\text{single}(\rho)$  which computes a formula whose denotation includes only singleton sequences of focused tree nodes contained in  $\llbracket \rho \rrbracket$ . That is,  $(\varphi_{\text{root}}, u) <: \rho$  if and only if  $\llbracket \varphi_{\text{root}} \wedge \text{form}(u) \wedge \neg \text{single}(\rho) \rrbracket = \emptyset$  which can be tested in  $2^{O(|u|+|\rho|)}$  time by the decision procedure in [28].

As for an if-expression  $\text{if } \text{nempty}(e_1) \text{ then } e_2 \text{ else } e_3$ , we first check if the condition expression  $e_1$  always reduces to a non-empty sequence or an empty sequence, regardless of the input trees. In other words, we first compute  $\mathcal{S} \leftarrow e_1 : (\top, \text{AnyElt})^+$  and  $\mathcal{S}' \leftarrow e_1 : ()$ . Then, if  $\mathcal{S}'$  (or  $\mathcal{S}$ ) is unsatisfiable, then we use the rule I-IFNONEMPTY (or the rule I-IFEMPTY). If both  $\mathcal{S}$  and  $\mathcal{S}'$  are satisfiable, that is,  $e_1$  can reduce to both a non-empty sequence and an empty sequence, depending on the input trees, then we use the rule I-IFANY. It simply assumes that  $e_1$  reduces to any sequence and infers a constraint  $\mathcal{S}_i$  from each subexpression  $e_i$ . If all of  $\mathcal{S}_1, \mathcal{S}_2$ , and  $\mathcal{S}_3$  are satisfiable, then the if-expression has the specified output type  $\rho$ .

For a let-binding  $\text{let } \$\bar{v} := e_1 \text{ return } e_2$ , the rule I-LET first infers a constraint  $\mathcal{S}_2$  from  $e_2$  and the given output type  $\rho$ . Then, for each constraint-set  $C \in \mathcal{S}_2$  such that  $\llbracket C(\$ \bar{v}) \rrbracket \neq \emptyset$ , we infer a constraint  $\mathcal{S}_1$  from  $e_1$  and  $C(\$ \bar{v})$ . Note that if  $\llbracket C(\$ \bar{v}) \rrbracket = \emptyset$ , then  $C$  is unsatisfiable. In order for the whole let-expression to have type  $\rho$ , both  $\mathcal{S}_1$  and  $C \setminus \$\bar{v}$  should be satisfiable, i.e.,  $\mathcal{S}_1 \sqcap \{C \setminus \$\bar{v}\}$ , where  $C \setminus \$\bar{v}$  removes the constraint for  $\$ \bar{v}$  from  $C$  because it is bound only in  $e_2$ .

For a sequence concatenation, the rule I-SEQ divides the output type  $\rho$  into two parts using an auxiliary function  $\text{split}()$ , defined in Fig. 19. Then, we infer an input constraint for each case in  $\text{split}(\rho)$ , and returns a union of all inferred constraints as a final result. Note that for any  $\rho$ , if  $(\rho_1, \rho_2) \in \text{split}(\rho)$  then  $\rho_1, \rho_2 <: \rho$ .

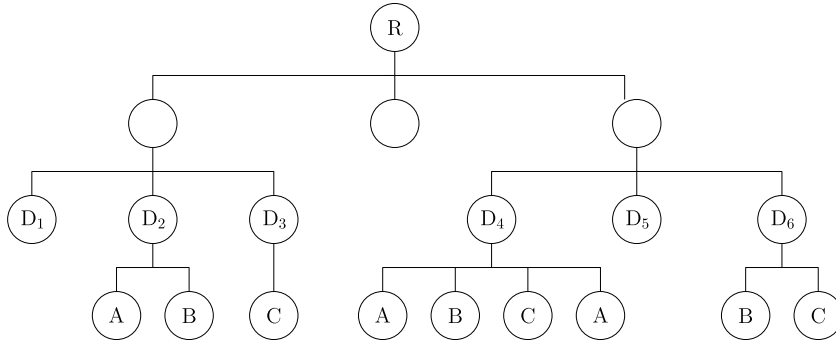
Finally, let us consider for-loop expressions, which are the most challenging with respect to defining precise inference rules. Indeed, they are the main source of the approximation introduced in our backward type inference. To illustrate, consider the following expression:

for  $\$v$  in  $\$doc/\$desc :: D$  return  $\$v/\$child :: *$

where  $\$doc$  is bound to an input tree. If  $\$doc/\$desc :: D$  reduces to  $[f_1, \dots, f_n]$  for some  $n$ , then the whole expression reduces to  $f_1/\$child :: *, \dots, f_n/\$child :: *$ . Suppose that an output type  $\rho_0$  is given as follows:

$$\rho_0 \equiv \langle A / \rangle \langle B / \rangle \langle C / \rangle \langle A / \rangle \langle B / \rangle \langle C / \rangle \langle A / \rangle \langle B / \rangle \langle C / \rangle$$

where for simplicity we use  $\langle A / \rangle$  to denote a formula type  $(\top, \text{element } A \{ () \})$  and juxtaposition to denote a sequence concatenation. To infer the exact type of  $\$doc$ , we need to infer the exact type of each  $f_i$ . Since the output sequence type



**Fig. 20.** An example XML tree bound to the variable  $\$doc$ : a for-loop expression  $\text{for } \$v \text{ in } \$doc/\text{desc}::D \text{ return } \$v/\text{child}::*$  reduces to a sequence of focused tree nodes  $[A[\epsilon], B[\epsilon], C[\epsilon], A[\epsilon], B[\epsilon], C[\epsilon], A[\epsilon], B[\epsilon], C[\epsilon]]$  where we omit the context of each node for simplicity. In addition, the for-loop expression may have type  $\langle A/\rangle\langle B/\rangle\langle C/\rangle^+$ , or less precisely  $\langle A/\rangle|\langle B/\rangle|\langle C/\rangle^*$ .

$\rho_0$  is finite, it suffices to compute all possible (weak) compositions of  $\rho_0$ , infer an input constraint for each composition, and take a union of all inferred constraints [33]. We say that  $(\rho_1, \dots, \rho_n)$  is a composition of  $\rho$  if  $(\rho_1, \dots, \rho_n) <: \rho$ . A composition is said to be *weak* if it contains an empty sequence type as an element. For example, we can infer the exact type of the example input tree given in Fig. 20 from the following weak composition:

$$(\ () \ ; \ \langle A/\rangle\langle B/\rangle \ ; \ \langle C/\rangle \ ; \ \langle A/\rangle\langle B/\rangle\langle C/\rangle\langle A/\rangle \ ; \ (\ ) \ ; \ \langle B/\rangle\langle C/\rangle)$$

For this input tree,  $\$doc/\text{desc}::D$  reduces to  $[f_1, \dots, f_6]$  where each  $f_i$  is a focused tree rooted at the node labeled  $D_i$  (here, the subscript  $i$  is not part of the node label; it is used solely to distinguish the nodes with the same label). Note that the composition above consists of the exact type of each  $f_i/\text{child}::*$ , given the output type  $\rho_0$ . Assume that  $\rho_i$  is inferred from  $f_i/\text{child}::*$  and the corresponding type in the composition given above. Then, our backward type inference infers the following exact input type for the whole for-loop expression (with some simplification):

$$\left( \text{fstDescFoll}(\rho_1 \wedge \text{fstDescFoll}(\rho_2 \wedge \text{fstDescFoll}(\rho_3 \wedge \text{fstDescFoll}(\rho_4 \wedge \text{fstDescFoll}(\rho_5 \wedge \text{fstDescFoll}(\rho_6 \wedge \text{noNextUpTo}(D, \alpha), D), D), D), D), D), D) \wedge \text{noWhereElse}(\alpha), \text{AnyElt} \right)$$

which states that the input node has six descendants with label  $D$  and the first descendant in document order is of type  $\rho_1$ , the second is of type  $\rho_2$ , and so on.

The situation becomes more complex if we consider repetition types. In the presence of repetition operators, the number of possible compositions is infinite in general.<sup>2</sup> To illustrate, consider an output type  $\rho'_0$  defined as follows:

$$\rho'_0 \equiv \langle A/\rangle\langle B/\rangle\langle C/\rangle^+$$

To infer the exact type of the input tree in Fig. 20 again, we need to unfold  $\rho'_0$  three times to obtain  $\rho_0$ , which is defined above, and compute its weak compositions. The problem is that in general we do not know statically how many times we need to unfold the given output repetition type to infer the exact input type for a for-loop expression. One possible solution is to unfold repetition types up to some arbitrary fixed number of times, giving up exact typing. Then the problem is to find such an unfolding number that allows practical and precise type inference.

In this paper, we adopt a simpler but more approximate approach. More precisely, we do not analyze output types of the form  $\rho^+$  or  $\rho^*$  across the boundary of  $\rho$ , that is, we do not unfold repetition types. We also do not compute compositions even for simple output sequence types containing no repetition operators. Instead, given a for-loop expression, we simply consider only those cases where each execution of the `return` expression evaluates to a sequence of focused trees whose type is a subtype of the given output type. Consequently, the input tree in Fig. 20 is not accepted by our type system if the output type is given as  $\langle A/\rangle\langle B/\rangle\langle C/\rangle^+$ . Our system accepts only those input trees whose descendants labeled  $D$  have no child or children of type  $\langle A/\rangle\langle B/\rangle\langle C/\rangle^+$ . Nevertheless, the input tree in Fig. 20 is accepted if a more general output type is given such as  $\langle A/\rangle|\langle B/\rangle|\langle C/\rangle^*$ .

Our approximation is similar in spirit to the approximation used in forward type inference systems [8,17]. To illustrate, consider a regular tree type  $T$  defined recursively as follows, where we use  $\langle A/\rangle$  to denote a unit type element  $A \{ \ () \}$ :

$$T \equiv \text{element } C \{ \langle A/\rangle, T, \langle B/\rangle \} | \ ()$$

<sup>2</sup> Even in the absence of repetition operators, the number of naive weak compositions is infinite. For the purpose of typing for-loop expressions, however, it suffices to consider only strong compositions and assume that any number of empty sequence types may exist between every two adjacent single formula types used in each strong composition. For example, the following compositions can all be treated equally for type inference:  $\langle A/\rangle; \langle B/\rangle$ ,  $\langle A/\rangle; \ ()$ ;  $\langle B/\rangle$ , and so on.

In most forward type inference systems, if  $f$  has type  $T$ , then the type of  $f/\text{desc-or-self}::*$  is deduced as  $(T \mid \langle A/\rangle \mid \langle B/\rangle)^*$ . The exact type, however, is the union of  $(T, \langle A/\rangle)^n, (\langle B/\rangle)^n$  for  $n \in \mathbb{N}$ , which is not regular. In general, in forward type inference systems, a sequence type  $(u_1, \dots, u_n)^*$  is often approximated into a less precise type  $(u_1 \mid \dots \mid u_n)^*$ , losing the information on the order of elements. Similarly, in our backward type inference system, in order to accept more input trees, the output type should be given such that the order of elements does not matter.

To further clarify the consequence of our approximation, consider the following examples:

- (1) for  $\$v$  in  $\langle A/\rangle, \langle A/\rangle$  return  $\$v: \langle A/\rangle \langle A/\rangle$
- (2) for  $\$v$  in  $\langle A/\rangle \langle B/\rangle, \langle C/\rangle, \langle D/\rangle \langle A/\rangle$  return  $\$v/\text{child}::*: \langle B/\rangle \langle C/\rangle \langle D/\rangle$

Our backward type inference rejects the first example while accepts the second one. More specifically, the first example is rejected because the for-loop variable  $\$v$  can be bound only to a singleton sequence, i.e.,  $\$v$  cannot have type  $\langle A/\rangle \langle A/\rangle$ . (If the composition-based approach is used, then the first example is also accepted.) The second example is accepted because  $\$v/\text{child}::*$  can have type  $\langle B/\rangle \langle C/\rangle \langle D/\rangle$ . More specifically, our inference system infers an input type element  $* \{ \text{AnyElt}^*, \langle B/\rangle, \text{AnyElt}^*, \langle C/\rangle, \text{AnyElt}^*, \langle D/\rangle, \text{AnyElt}^* \}$  for  $\$v$  using the rule I-AXIS (with some simplification) which matches with  $\langle A/\rangle \langle B/\rangle, \langle C/\rangle, \langle D/\rangle \langle A/\rangle$ . However, the following similar examples are not accepted because the result of each execution of the return expression does not have the specified output type  $\langle B/\rangle \langle C/\rangle \langle D/\rangle$  or  $(\langle B/\rangle \langle C/\rangle \langle D/\rangle)^+$ .

- (3) for  $\$v$  in  $\langle A/\rangle \langle B/\rangle \langle A/\rangle, \langle A/\rangle, \langle A/\rangle \langle C/\rangle \langle D/\rangle \langle A/\rangle$  return  $\$v/\text{child}::*: \langle B/\rangle \langle C/\rangle \langle D/\rangle$
- (4) for  $\$v$  in  $\langle A/\rangle \langle B/\rangle \langle A/\rangle, \langle A/\rangle, \langle A/\rangle \langle C/\rangle \langle D/\rangle \langle A/\rangle$  return  $\$v/\text{child}::*: (\langle B/\rangle \langle C/\rangle \langle D/\rangle)^+$

Still, the first and third examples are accepted if they are given a more general output type as follows:

- (5) for  $\$v$  in  $\langle A/\rangle, \langle A/\rangle$  return  $\$v: \langle A/\rangle^+$
- (6) for  $\$v$  in  $\langle A/\rangle \langle B/\rangle \langle A/\rangle, \langle A/\rangle, \langle A/\rangle \langle C/\rangle \langle D/\rangle \langle A/\rangle$  return  $\$v/\text{child}::*: (\langle B/\rangle \mid \langle C/\rangle \mid \langle D/\rangle)^+$

We will explain how these examples are accepted by our type inference system after discussing the inference rules for for-loop expressions below.

To infer an input type for for-loop expressions, we use three rules. First, given an expression for  $\$v$  in  $e_1$  return  $e_2$ , if the output type is  $()$ , then no matter how many times we evaluate  $e_2$  with different bindings for  $\$v$ , it must reduce to  $\epsilon$ . Therefore, in the rule I-FOREMPTY, we infer a constraint  $\mathcal{S}_1$  from  $e_1$  with the output type  $C(\$v)^*$  where  $C$  is a constraint-set inferred by analyzing  $e_2$  with  $()$ . Note that  $C(\$v)^*$  is zero or more repetitions of a single formula type  $C(\$v)$  that makes  $e_2$  reduce to an empty sequence. In addition, we also consider the case where  $e_1$  reduces to an empty sequence  $\epsilon$  since if  $e_1$  reduces to  $\epsilon$ , then the whole for-loop expression also reduces to  $\epsilon$ . The remaining two rules cover the cases where the output type  $\rho$  is not  $()$ . In particular, the rules I-FOR and I-FORNULLABLE cover the cases where the output type  $\rho$  is not nullable and nullable, respectively.

The rule I-FOR combines the following three cases. First, the set  $\mathcal{S}'$  of constraint-sets, inferred from  $e_2$  and  $()$ , is unsatisfiable, i.e.,  $\mathcal{S}' = \emptyset$ . This means that  $e_2$  never reduces to an empty sequence regardless of the value of  $\$v$ . Second, every constraint-set inferred from  $e_2$  and  $\rho$  is incompatible with every constraint-set inferred from  $e_2$  and  $()$ , i.e.,  $C \setminus_{\$v} \sqcap C' \setminus_{\$v} = \{\perp\}$  for every  $(C, C') \in \mathcal{S} \times \mathcal{S}'$ . Finally,  $e_2$  can reduce to both an empty and a non-empty sequence depending on the value of  $\$v$ . The first and second cases are treated equally and thus we explain them first. For these cases, the set  $\mathcal{S}'$  is empty, and we use only the constraint-sets inferred from  $e_2$  and  $\rho$ , i.e., those in  $\mathcal{S}$ . To illustrate, suppose that  $e_1$  reduces to  $f_1, \dots, f_n$ . Then, we accept only the cases where  $\llbracket e_2 \rrbracket_{\eta, \$v \mapsto f_i}$  has type  $\rho$  for all  $i$ 's and for some substitution  $\eta$ . In other words, for some  $i$ , if  $\llbracket e_2 \rrbracket_{\eta, \$v \mapsto f_i}$  has type  $\rho'$  which is not a subtype of  $\rho$ , then the given for-loop expression is rejected by our inference system. More specifically, if the output type  $\rho$  is not a repetition type, then the rule I-FOR requires  $e_1$  to reduce to a single focused tree of type  $C(\$v)$  where  $C$  is a constraint-set inferred from  $e_2$  and  $\rho$ . However, if  $e_1$  may reduce to a sequence of more than one focused tree node, say  $f_1, f_2$ , then it is rejected since the whole for-loop expression would reduce to a sequence  $\llbracket e_2 \rrbracket_{\eta, \$v \mapsto f_1}, \llbracket e_2 \rrbracket_{\eta, \$v \mapsto f_2}$  whose type is  $(\rho, \rho)$  which is not equivalent to  $\rho$ , that is, the output type of the for-loop expression. Now consider the case where the output type  $\rho$  is a repetition type of the form  $\rho_0^+$ . In this case,  $e_1$  may reduce to a sequence of any number of focused trees because no matter how many times we evaluate  $e_2$ , it would reduce to a sequence of type  $\rho_0^+$  and their concatenation is also of type  $\rho_0^+$  (i.e.,  $\rho_0^+, \dots, \rho_0^+ <: \rho_0^+$ ). Therefore, the output type of  $e_1$  is determined as  $C(\$v)^+$ , that is, we require  $e_1$  to reduce to a sequence of nodes  $f_1, \dots, f_n$  where each  $f_i$  is of type  $C(\$v)$ . Here  $C$  is again a constraint-set inferred from  $e_2$  and  $\rho$ . To distinguish repetition types from other types, we use an auxiliary function  $\text{Quant}()$ . Finally, the rule I-FOR also considers the case where  $e_2$  can reduce to both an empty and a non-empty sequence depending on the value of  $\$v$ . In this case, we also infer an input constraint from  $e_1$  with the output sequence type  $(C'(\$v)^*, C(\$v), C'(\$v)^*) \cdot \text{Quant}(\rho)$  where  $C'(\$v)$  and  $C(\$v)$  are the types that make  $e_2$  reduce to an empty and a non-empty sequence, respectively. The rule I-FORNULLABLE is the same as the rule I-FOR except that it also considers the case where  $e_1$  reduces to an empty sequence and  $\rho$  is nullable, as in the rule I-FOREMPTY.

Now let us consider the example (5) given above. For this example, we apply the rule I-FOR since the given output type is not nullable. Note that  $\emptyset \leftarrow \$v: ()$ . Next, a constraint-set  $\{\$v: \langle A/\rangle\}$  is inferred by the rule I-FVAR with  $\$v$  and  $\langle A/\rangle^+$ .

Then,  $\langle A/\rangle, \langle A/\rangle$  is matched against  $\langle A/\rangle^+$  which succeeds by the rules I-SEQ and I-ELEMENT. As for the example (6), we also apply the rule I-FOR. Note that  $\$/\text{child}::*$  can be of type  $()$ , i.e., when  $\$/$  is of type  $\neg\langle 1\rangle\top$  which means that there are no child nodes. Next, from  $\$/\text{child}::*$  and  $(\langle B/\rangle|\langle C/\rangle|\langle D/\rangle)^+$ , a constraint-set  $\{\$/\text{child}::*\}$  is inferred by the rule I-AXIS where  $T$  is defined as follows (with some simplification):

$$T \equiv \text{element} * \{(\langle B/\rangle|\langle C/\rangle|\langle D/\rangle)^+\}$$

Then,  $\langle A/\rangle\{\langle B/\rangle\}\langle A/\rangle, \langle A/\rangle, \langle A/\rangle\{\langle C/\rangle\langle D/\rangle\}\langle A/\rangle$  is matched against  $((\neg\langle 1\rangle\top)^*, T, (\neg\langle 1\rangle\top)^*)^+$  which again succeeds by the rules I-SEQ and I-ELEMENT.

## 5.2. Complexity

### 5.2.1. Complexity for XPath axes

We first analyze the complexity of our backward type inference system for XPath axes. To this end, we first define the length  $\text{len}(\rho)$  and the size  $|\rho|$  of a formula-enriched sequence type  $\rho$ :

$$\begin{array}{ll} \text{len}((\varphi, u)) &= 1 & |(\varphi, u)| &= |\varphi| + |u| \\ \text{len}() &= 1 & |()| &= 1 \\ \text{len}(\rho_1, \rho_2) &= \text{len}(\rho_1) + \text{len}(\rho_2) + 1 & |\rho_1, \rho_2| &= |\rho_1| + |\rho_2| + 1 \\ \text{len}(\rho_1 | \rho_2) &= \text{len}(\rho_1) + \text{len}(\rho_2) + 1 & |\rho_1 | \rho_2| &= |\rho_1| + |\rho_2| + 1 \\ \text{len}(\rho^+) &= \text{len}(\rho) + 1 & |\rho^+| &= |\rho| + 1 \end{array}$$

The size  $|\varphi|$  of a formula  $\varphi$  and the length  $\text{len}(\tau)$  and the size  $|\tau|$  of a regular tree type  $\tau$  are also defined as usual. In particular, in the analysis below, we mean by  $|\tau|$  the size of the classical binary representation of  $\tau$  [2].

**Lemma 5.2.** *The time complexity of computing an application of each auxiliary function introduced in Section 4 is as follows.*

- $\text{nullable}(\rho)$  can be computed in  $O(\text{len}(\rho))$  time.
- $\text{child-type}(\rho)$  can be computed in  $O(\text{len}(\rho))$  time.
- $\text{parent-type}(\rho)$  can be computed in  $O(\text{len}(\rho))$  time.
- $\text{desc-type}(\rho)$  can be computed in  $O(|\rho|)$  time.

$\text{child-type}(\rho)$  in Fig. 4 is defined only when the argument  $\rho$  is of the form  $(\varphi_1, u_1) | \dots | (\varphi_n, u_n)$  where  $u_i = \text{element } n_i \{ \tau_i \}$ , and its precise complexity is indeed  $O(\text{len}(\rho) \times \max \text{len}(\tau_i))$ . We consider  $\max \text{len}(\tau_i)$  as a constant and omit it in the above analysis since it is usually small and does not affect the overall complexity of our inference system.

Among the functions listed in Lemma 5.2, only  $\text{nullable}()$  may be called many times during the inference. More precisely, when the output type is  $(\varphi_1, u_1), \dots, (\varphi_n, u_n)$ , the naive cumulative cost of calling  $\text{nullable}()$  is in total  $O(n^2)$ . With additional space, however, if we memoize the result of  $\text{nullable}()$  on each subterm of the output type  $\rho$  when it is called for the first time, the cumulative cost is still  $O(\text{len}(\rho))$ .

**Lemma 5.3.** *Given an output type  $\rho$ , an input type for an XPath axis is inferred in  $O(|\rho|)$  time.*

**Proof.** Easy from the fact that we analyze the structure of the output type, with an empty type  $()$  and a pair type  $(\varphi, u)$  as base cases, and the cumulative cost of using auxiliary functions during the inference is  $O(|\rho|)$ .  $\square$

To analyze the size of the inferred input type, below we assume that we use an optimization technique such as hashing to represent types and formulas, i.e., to share the same subterms. Otherwise, in the input type, some formula may be duplicated an exponential number of times in terms of the length of the output type, e.g., when the output type is of the form  $(\rho_1 | \rho_2), \dots, (\rho_{n-1} | \rho_n)$ . Note that in the rule AXIS-OR, with a naive representation of formulas, the with parameter  $\psi$  may be duplicated in the inferred input formula  $\varphi_1 \vee \varphi_2$ : one in  $\varphi_1$ , the other in  $\varphi_2$ .

**Lemma 5.4.** *Assume  $\varphi \leftarrow \text{axis}::n, \rho$  with  $\psi$ . Then the size of  $\varphi$  is  $O(|\rho| + |\psi|)$ .*

**Proof.** By induction on a derivation of  $\varphi \leftarrow \text{axis}::n, \rho$  with  $\psi$ . In the proof, we use the fact that all the auxiliary definitions used in Fig. 16, which take a formula  $\chi$  as argument, return another formula of size  $O(|\chi|)$ . The proof also relies on that  $\text{form}(u)$  has the same size as the classical binary representation of the regular tree type  $u$  [19].  $\square$

**Lemma 5.5.** *Given an output type  $\rho$ , the size of the inferred input type for an XPath axis is  $O(|\rho|)$ .*



**Proof.** The cases for the axes except `self`, `parent`, `child`, and `anc` are easily proved by Lemma 5.4. The case for `self` is proved by structural induction on the output type  $\rho$ . The cases for `parent`, `child`, and `anc` are proved by the fact that the size of `child-type`( $\rho$ ), `parent-type`( $\rho$ ), and `desc-type`( $\rho$ ) is  $O(|\rho|)$  with optimized representations of types.  $\square$

**Corollary 5.6.** Given an output type  $\rho$  and an XPath axis, we can check in  $2^{O(|\rho|)}$  time if there exists some tree that when applied to the axis, returns a sequence of nodes of type  $\rho$ , by testing the satisfiability of the inferred input type using the decision procedure in [28].

Precisely, if  $\rho' \leftarrow \text{axis} : n, \rho$ , then  $\rho'$  is of the form  $(\varphi_1, u_1) \mid \dots \mid (\varphi_n, u_n)$  where  $\varphi_i <: u_i$ , and thus it suffices to check the satisfiability of each  $\varphi_i$  in the inferred input type.

### 5.2.2. Complexity for the XQuery core

Now we analyze the complexity of our backward type inference system for the XQuery core. We define the size  $|C|$  of  $C$  and the size  $|\mathcal{S}|$  of  $\mathcal{S}$  as the number of bindings in  $C$  and the number of constraint-sets in  $\mathcal{S}$ , respectively. In particular,  $|C|$  does not take into consideration the size of the types included in  $C$  and similarly for  $|\mathcal{S}|$ . Then,  $|C_1 \sqcap C_2| \leq |C_1| + |C_2|$ ,  $|\mathcal{S}_1 \sqcap \mathcal{S}_2| \leq |\mathcal{S}_1| \times |\mathcal{S}_2|$ , and  $|\mathcal{S}_1 \sqcup \mathcal{S}_2| \leq |\mathcal{S}_1| + |\mathcal{S}_2|$ . The size  $|e|$  of an XQuery expression  $e$  is inductively defined as usual, e.g., see Definition 8.1 in [8].

**Lemma 5.7.** Suppose  $\mathcal{S} \leftarrow e : \rho$ . Then the maximum size, denoted by  $T(e, \rho)$ , of a largest type appearing in  $\mathcal{S}$  is  $O(2^{|e|}|\rho|)$ , i.e., single exponential in terms of the size of the given expression  $e$ .

**Proof.** By solving the following set of recursive equations, which are derived from the inference rules:

$$\begin{aligned}
T(e, \rho_1 \mid \rho_2) &= \max_i T(e, \rho_i) \\
T(e, \rho_1 \wedge \rho_2) &= T(e, \rho_1) + T(e, \rho_2) + 1 \\
T((e_1, e_2), \rho) &= \max_{(\rho_1, \rho_2) \in \text{split}(\rho)} (T(e_1, \rho_1) + T(e_2, \rho_2) + 1) \\
T((\langle \sigma \rangle \{e\} \langle / \sigma \rangle : \text{element } n \{ \tau \}), \rho) &= T(e, \text{form-enriched}(\tau)) \\
T(\text{if } \text{nempty}(e_1) \text{ then } e_2 \text{ else } e_3, \rho) &= T(e_1, (\top, \text{AnyElt})^*) + T(e_2, \rho) + T(e_3, \rho) + 2 \\
T(\text{let } \$\bar{v} := e_1 \text{ return } e_2, \rho) &= T(e_2, \rho) + T(e_1, T(e_2, \rho)) + 1 \\
T(\text{for } \$v \text{ in } e_1 \text{ return } e_2, \rho) &= T(e_2, \rho) + T(e_2, ()) + T(e_1, (T(e_2, \rho) + 2T(e_2, ()) + 5)) + 2 \\
T(e, \rho) &= O(|\rho|) \quad (\text{otherwise})
\end{aligned}$$

where we use a type and its size interchangeably as the second argument to  $T(-, -)$ .  $\square$

**Lemma 5.8.** Suppose  $\mathcal{S} \leftarrow e : \rho$ . Then the maximum size, denoted by  $N(e, \rho)$ , of  $\mathcal{S}$  is  $O(2^{2^{|e|}|\rho|})$ , i.e., double exponential in terms of the size of the given expression  $e$ .

**Proof.** By solving the following set of recursive equations, which are derived from the inference rules. We use the result from Lemma 5.7.

$$\begin{aligned}
N(e, \rho_1 \mid \rho_2) &= N(e, \rho_1) + N(e, \rho_2) \\
N(e, \rho_1 \wedge \rho_2) &= N(e, \rho_1) \times N(e, \rho_2) \\
N((e_1, e_2), \rho) &= |\text{split}(\rho)| \times \max_{(\rho_1, \rho_2) \in \text{split}(\rho)} (N(e_1, \rho_1) \times N(e_2, \rho_2)) \\
N((\langle \sigma \rangle \{e\} \langle / \sigma \rangle : \text{element } n \{ \tau \}), \rho) &= N(e, \text{form-enriched}(\tau)) \\
N(\text{if } \text{nempty}(e_1) \text{ then } e_2 \text{ else } e_3, \rho) &= N(e_1, (\top, \text{AnyElt})^*) \times N(e_2, \rho) \times N(e_3, \rho) \\
N(\text{let } \$\bar{v} := e_1 \text{ return } e_2, \rho) &= N(e_2, \rho) \times N(e_1, T(e_2, \rho)) \\
N(\text{for } \$v \text{ in } e_1 \text{ return } e_2, \rho) &= N(e_2, \rho) \times N(e_2, ()) \times N(e_1, (T(e_2, \rho) + 2T(e_2, ()) + 5)) + \\
&\quad N(e_2, \rho) \times N(e_1, T(e_2, \rho)) \\
N(e, \rho) &= O(1) \quad (\text{otherwise})
\end{aligned}$$

In the above equations, we use a type and its size interchangeably as the second argument to  $N(-, -)$ .  $\square$

**Lemma 5.9.** Suppose  $\mathcal{S} \leftarrow e : \rho$ . Then  $\mathcal{S}$  is computed in  $2^{O(2^{|e|}|\rho|)}$  time in the worst case.

**Proof.** Let  $I(e, \rho)$  denote the complexity of deducing a set of constraint-sets from  $e$  and  $\rho$  using our inference system. We obtain the complexity by solving the following set of recursive equations, which are derived from the inference rules. We use the result from Lemmas 5.7 and 5.8. We also use a type and its size interchangeably as the second argument to  $I(-, -)$ .

$$\begin{aligned}
I(\epsilon, \rho) &= I(\overline{\$v}, \rho) = 1 \\
I(e, \rho_1 \mid \rho_2) &= I(e, \rho_1 \wedge \rho_2) = 1 + I(e, \rho_1) + I(e, \rho_2) \\
I(\$v, \rho) &= I(\$v/\text{axis}:n, \rho) = O(|\rho|) \\
I((e_1, e_2), \rho) &= 1 + |\text{split}(\rho)| \times \max_{(\rho_1, \rho_2) \in \text{split}(\rho)} (I(e_1, \rho_1) + I(e_2, \rho_2)) \\
I(\langle \langle \sigma \rangle \{e\} \langle / \sigma \rangle : \text{element } n \{ \tau \} \rangle, \rho) &= 2 + I(e, \text{form-enriched}(\tau)) + 2^{O(|u|+|\rho|)} \\
I(\text{if nempty}(e_1) \text{ then } e_2 \text{ else } e_3, \rho) &= 1 + I(e_1, (\top, \text{AnyElt})^*) + I(e_2, \rho) + I(e_3, \rho) \\
I(\text{let } \overline{\$v} := e_1 \text{ return } e_2, \rho) &= 1 + I(e_2, \rho) + N(e_2, \rho) \times I(e_1, T(e_2, \rho)) + N(e_2, \rho) \times 2^{O(T(e_2, \rho))} \\
I(\text{for } \$v \text{ in } e_1 \text{ return } e_2, \rho) &= 1 + I(e_2, \rho) + I(e_2, ()) + O(\text{len}(\rho)) + I(e_1, ()) + \\
&\quad N(e_2, \rho) \times N(e_2, ()) \times I(e_1, (T(e_2, \rho) + 2T(e_2, ()) + 5)) + \\
&\quad N(e_2, \rho) \times I(e_1, T(e_2, \rho)) + N(e_2, \rho) \times 2^{O(T(e_2, \rho))} + \\
&\quad N(e_2, ()) \times 2^{O(T(e_2, ()))}
\end{aligned}$$

In the above equations, the case of the element construction includes the complexity for the subtype check  $u <: \rho$ . The cases of let-expressions and for-loop expressions include the complexity of satisfiability checks for the inferred type for the bound variable, e.g.,  $C(\$var)$ .  $\square$

Lastly, we state the worst-case time complexity of our backward type inference for the XQuery core.

**Theorem 5.10** (Complexity of type inference). *Assume we are given an XQuery expression  $e$  and its output type  $\rho$ . Then a set of solvable constraint-sets is computed in  $|e| \cdot 2^{O(2^{(|e|+1)|\rho|})}$  time by our inference system. That is, the overall cost is double exponential in terms of the size of the given expression  $e$ .*

**Proof.** Suppose  $\mathcal{S} \leftarrow e : \rho$ . We obtain  $\mathcal{S}$  in  $2^{O(2^{|\rho|})}$  time by Lemma 5.9. The size of  $\mathcal{S}$  is  $O(2^{2^{|\rho|}})$  by Lemma 5.8. The size of any constraint-set  $C$  in  $\mathcal{S}$  is the number of free variables in  $e$ , which is at most  $|e|$ . Since the size of the largest type in  $\mathcal{S}$  is  $O(2^{|\rho|})$  by Lemma 5.7, for each constraint-set  $C$  in  $\mathcal{S}$ , its satisfiability can be tested in  $|e| \cdot 2^{O(2^{|\rho|})}$  time by the decision procedure in [28]. Overall, the complexity of our inference system is  $2^{O(2^{|\rho|})} + O(2^{2^{|\rho|}}) \times |e| \cdot 2^{O(2^{|\rho|})}$  which is simply  $|e| \cdot 2^{O(2^{(|e|+1)|\rho|})}$ .  $\square$

### 5.3. Soundness

Now we state the soundness property for our inference system. Below we use  $\vdash \eta : C$  to mean that if  $\$var \mapsto s \in \eta$ , then  $(\$var : \rho) \in C$  and  $s \in \llbracket \rho \rrbracket$ .

**Theorem 5.11** (Soundness). *Let  $e$  and  $\rho$  be an XQuery expression and its output type, respectively. Suppose  $\mathcal{S} \leftarrow e : \rho$ . Then for any  $C \in \mathcal{S}$  such that  $C \neq \{\perp\}$ , if  $\vdash \eta : C$  and  $\llbracket e \rrbracket_\eta = s$ , then  $s \in \llbracket \rho \rrbracket$ .*

**Proof.** By induction on a derivation of  $\mathcal{S} \leftarrow e : \rho$ . Here we only show the case for the rule I-FOR. Other cases are similarly proved. We have the following assumptions:

- (1)  $\sqcup_{\mathcal{S} \in \text{SUS}' \mathcal{S}} \mathcal{S} \leftarrow \text{for } \$v \text{ in } e_1 \text{ return } e_2 : \rho$
- (2)  $C_0 \in \sqcup_{\mathcal{S} \in \text{SUS}' \mathcal{S}}$  and  $\vdash \eta : C_0$
- (3)  $\llbracket \text{for } \$v \text{ in } e_1 \text{ return } e_2 \rrbracket_\eta = s$

Then, we need to prove  $s \in \llbracket \rho \rrbracket$ .

- (4) Let  $\sqcup_{\mathcal{S} \in \text{SUS}' \mathcal{S}}$  be  $\mathcal{S}_1 \sqcup \dots \sqcup \mathcal{S}_m$ .
- (5) Without loss of generality, let  $C_0 \in \mathcal{S}_i \in \text{S}'$ .

The case for  $\mathcal{S}_i \in \text{S}$  is similarly proved. From the premises of the rule I-FOR, we have

- (6)  $\mathcal{S} \leftarrow e_2 : \rho$
- (7)  $\mathcal{S}' \leftarrow e_2 : ()$
- (8)  $C \in \mathcal{S}$  and  $C' \in \mathcal{S}'$
- (9)  $\mathcal{S}'' \leftarrow e_1 : (C'(\$v)^*, C(\$v), C'(\$v)^*). \text{Quant}(\rho)$
- (10)  $\mathcal{S}_i = \mathcal{S}'' \sqcap \{C \setminus_{\$v} \sqcap C' \setminus_{\$v}\}$

From (3), we have

- (11)  $\llbracket e_1 \rrbracket_\eta = f_1, \dots, f_n$

$$(12) s = \prod_{f_1, \dots, f_n} \llbracket e_2 \rrbracket_{\eta, \$v \mapsto f_i}$$

From (2), (5) and (10),

$$(13) \exists C'_0 \in \mathcal{S}'' \text{ such that } C_0 = C'_0 \sqcap C \setminus \$v \sqcap C' \setminus \$v \text{ and } \vdash \eta : C'_0.$$

By induction hypothesis on (9) with (11) and (13), we have

$$(14) f_1, \dots, f_n \in \llbracket (C'(\$v)^*, C(\$v), C'(\$v)^*).Quant(\rho) \rrbracket.$$

Assume  $Quant(\rho) = 1$ . The case where  $Quant(\rho) = +$  is similarly proved using the following property:  $\rho^+, \dots, \rho^+ <: \rho^+$ . Then, there exists  $j$  such that

$$(15) f_1, \dots, f_{j-1} \in \llbracket C'(\$v)^* \rrbracket \text{ and thus } f_k \in \llbracket C'(\$v) \rrbracket \text{ where } k = 1, \dots, j-1$$

$$(16) f_j \in \llbracket C(\$v) \rrbracket$$

$$(17) f_{j+1}, \dots, f_n \in \llbracket C'(\$v)^* \rrbracket \text{ and thus } f_k \in \llbracket C'(\$v) \rrbracket \text{ where } k = j+1, \dots, n$$

From (2), (5) and (10), we have  $\vdash \eta : C \setminus \$v$  and  $\vdash \eta : C' \setminus \$v$ . Together with (15)–(17), we have

$$(18) \vdash \eta, \$v \mapsto f_j : C$$

$$(19) \vdash \eta, \$v \mapsto f_k : C' \text{ where } k = 1, \dots, j-1, j+1, \dots, n$$

By induction hypothesis on (6) and (7) with (18) and (19), respectively, we have

$$(20) \llbracket e_2 \rrbracket_{\eta, \$v \mapsto f_j} \in \llbracket \rho \rrbracket$$

$$(21) \llbracket e_2 \rrbracket_{\eta, \$v \mapsto f_k} \in \llbracket () \rrbracket \text{ where } k = 1, \dots, j-1, j+1, \dots, n$$

From (20) and (21), we have  $s = \prod_{f_1, \dots, f_n} \llbracket e_2 \rrbracket_{\eta, \$v \mapsto f_i} \in \llbracket \rho \rrbracket$  as desired.  $\square$

Unlike the type inference for XPath axes, the type inference for the XQuery core is only sound and not complete, mainly because of the approximation introduced for for-loop expressions. From the soundness and the decidability of the inference system, we deduce a sound typechecking algorithm as a corollary.

**Corollary 5.12** (Typechecking). *Let  $e$  be an XQuery expression with the only free variable  $\$doc$ , which denotes an input document. Let  $\rho_i$  be an input type (the type for  $\$doc$ ) and  $\rho_o$  an output type. Then there exists a typechecking algorithm that says yes if  $\mathcal{S} \leftarrow e : \rho_o$  and  $\exists C \in \mathcal{S}$  such that  $C \neq \{\perp\}$  and  $\rho_i <: C(\$doc)$ . Combined with the soundness property in Theorem 5.11, for any tree node  $t \in \llbracket \rho_i \rrbracket$ , if  $\llbracket e(t) \rrbracket = s$ , then  $s \in \llbracket \rho \rrbracket$  is guaranteed.*

In the corollary above, the input type  $\rho_i$  should be of the form  $(\varphi_1, u_1) \mid \dots \mid (\varphi_n, u_n)$  or simply  $u$  since it is the type for XML documents, but not for arbitrary XQuery expressions. Indeed, the inferred type  $C(\$doc)$  is also of the form  $(\varphi'_1, u'_1) \mid \dots \mid (\varphi'_m, u'_m)$ . Moreover, the use of the variable  $\$doc$  has no particular implication; it suffices to have a name of some element that is considered as the root in the input type, e.g.,  $e$  should be of the form  $\text{let } \$doc := /self:: * \text{ return } e'$ . To typecheck a given expression  $e$  with the input type  $\rho_i$  and the output type  $\rho_o$ , we first infer a constraint-set  $C$  from  $e$  and  $\rho_o$  using our backward type inference, and then simply check the inclusion relation between  $\rho_i$  and the inferred type  $C(\$doc)$ .

**Theorem 5.13** (Complexity of typechecking). *Let  $e$  be an XQuery expression with only one free variable  $\$doc$ . Then, given an input type  $\rho_i$  of the form  $(\varphi_1, u_1) \mid \dots \mid (\varphi_n, u_n)$  and an output type  $\rho_o$ ,  $e$  can be typechecked in  $2^{O(2^{|\rho_i|+|\rho_o|})}$  time. That is, the complexity of typechecking is double exponential in terms of the size of the given expression  $e$ .*

**Proof.** Suppose  $\mathcal{S} \leftarrow e : \rho_o$  where  $\mathcal{S}$  is a set of solvable constraint-sets.  $\mathcal{S}$  can be computed in  $2^{O(2^{|\rho_i|+|\rho_o|})}$  time by Theorem 5.10. (Note that there is only one free variable in  $e$ .) For the purpose of typechecking, we need to check if there exists  $C \in \mathcal{S}$  such that  $\rho_i <: C(\$doc)$  holds. We note that  $\rho_i \equiv (\varphi_1, u_1) \mid \dots \mid (\varphi_n, u_n)$  can be translated into an equivalent formula  $\psi \equiv (\varphi_1 \wedge \text{form}(u_1)) \vee \dots \vee (\varphi_n \wedge \text{form}(u_n))$  of the same size in terms of big O notation and similarly  $C(\$doc)$  into  $\psi'$ . Then, the subtype check  $\rho_i <: C(\$doc)$  can be done by testing the satisfiability of  $\psi \wedge \neg \psi'$  using the decision procedure in [28]. Since the size of the largest type appearing in  $\mathcal{S}$  is  $O(2^{|\rho_o|})$  by Lemma 5.7, the satisfiability test can be done in  $2^{O(2^{|\rho_o|+|\rho_i|})}$  time in the worst case. Since the size of  $\mathcal{S}$  is bounded by  $O(2^{2^{|\rho_o|}})$  by Lemma 5.8, we can check if there exists  $C \in \mathcal{S}$  such that  $\rho_i <: C(\$doc)$  holds in  $O(2^{2^{|\rho_o|}}) \times 2^{O(2^{|\rho_o|+|\rho_i|})}$  time, which is simply  $2^{O(2^{|\rho_i|+|\rho_o|})}$ . Overall, the complexity of typechecking is double exponential in terms of  $|e|$ .  $\square$

This result easily extends to the cases where more free variables are used. In such cases, we simply need an input type for each free variable and the complexity of typechecking still remains the same.

## 6. Related work and discussion

### 6.1. Typechecking for XML transformations

The problem of typechecking XML transformations has been extensively studied since the introduction of XML. There are two major approaches, namely forward type inference and backward type inference. Given an expression  $e$  that transforms XML documents of type  $\rho_i$  into documents of type  $\rho_o$ , forward type inference first computes the image  $O$  of the input type  $\rho_i$  under the transformation  $e$ , i.e.,  $O := \{e(t) \mid t \in \rho_i\}$ , and then checks if  $O \subseteq \rho_o$ . This approach does not work even for simple top-down tree transducers since  $O$  can be beyond context-free tree languages and in this case checking  $O \subseteq \rho_o$  is undecidable [14]. In contrast, backward type inference computes the pre-image  $J$  of the complement of  $\rho_o$  under  $e$ , i.e.,  $J := \{t \mid e(t) \in \neg\rho_o\}$ , and then checks the emptiness of intersection of  $J$  and  $\rho_i$ . (Or alternatively, for deterministic transducers, it may compute the pre-image  $I$  of the output type  $\rho_o$  and then checks if  $\rho_i \subseteq I$ .) When types are modeled as regular tree languages, exact typechecking may be done in the form of backward type inference by using tree transducers as a model of XML transformations [10–12]. In contrast, in forward type inference, even for simple XML transformations, their image may not be regular, as illustrated in Section 5.1. Therefore, if the type system infers regular tree types for admissible outputs, which is often the case in the literature, then typechecking cannot be exact. Still, forward type inference is more intuitive than backward type inference, and thus many practical XML programming languages such as XQuery [1,3], XDuCE [24], and CDuce [34] build on forward type inference and instead introduce some approximation, i.e., some type-safe programs are rejected in these languages. For a more detailed, general survey of typechecking for XML transformations, we refer the reader to [35,36] and references therein. Below we discuss only closely related work on backward type inference and precise type systems for XPath and XQuery.

### 6.2. Inverse type inference

A problem of *inverse type inference*, which is another name of backward type inference, has been extensively investigated to develop an exact typechecking algorithm for XML transformations [9–14]. For example, Milo et al. [10] propose an exact inverse type inference algorithm for *k-pebble tree transducers*, which are finite state transducers that can mark nodes of the input tree using up to  $k$  different pebbles. Although we can model a broad range of XQuery expressions using  $k$ -pebble tree transducers, the complexity of typechecking is hyper-exponential, i.e., when using  $k$  pebbles, its complexity is  $O(h_{k+2}(n))$  with  $h_0(n) = n$  and  $h_{m+1}(n) = 2^{h_m(n)}$ .

Maneth et al. [12] also study the problem of exact inverse type inference for tree transformations using macro tree transducers (MTTs) [15], which can accumulate part of the input and copy it in the output. Their transformation language called TL uses monadic second-order logic (MSO) as a pattern language, which subsumes XPath without arithmetics and data value comparisons. By using MTTs and MSO, TL can be used to describe many real-world XML transformations. Their formalism, however, is based on finite automata and thus requires for implementation purposes a translation from MSO to a finite automaton which may introduce a non-elementary blow-up.

Perst and Seidl [11] extend MTTs with concatenation and propose macro forest transducers (MFTs) as a model of XML transformations. They develop an exact inverse type inference algorithm for MFTs and show that the complexity of typechecking is DEXPTIME-complete. Moreover, by combining with a translation from a downward navigational fragment of XQuery into MFTs [37], MFTs can be used as an intermediate language for a subset of XQuery. The translation, however, considers only XPath axes such as `child`, `desc`, and `following-sibling`, and a restricted form of for-loop expressions, i.e., in `for $v in  $e_1$  return  $e_2$` ,  $e_1$  must be a path expression.

In order to support backward axes, one may use tree-walking automata [38] as a pattern language. Indeed, a  $k$ -pebble tree transducer can be decomposed into a  $(k+1)$ -fold composition of tree-walking transducers [16]. Similarly, a TL program using MSO patterns can be compiled into a composition of an MTT and a macro tree-walking transducer, which can then be decomposed into a three-fold composition of (stay) MTTs [12]. Therefore, the complexity of typechecking a TL program is quadruple exponential even if we do not consider a possible blow-up in the translation of MSO patterns to finite automata. In this paper, we also study the problem of backward type inference, but develop a type inference system directly on the XQuery core. We present an exact backward type inference algorithm for XPath axes whose complexity is simple exponential. This result corresponds to the fact that the complexity of inverse type inference for tree-walking transducers is also exponential [16]. As for the XQuery core, instead of trying to develop yet another hyper-exponential algorithm, we introduce a sound approximation similar to the one used in forward type systems.

### 6.3. Precise type systems for XPath and XQuery

Typing XPath expressions has been a challenging topic. Most previous proposals for the XQuery static type system, including the one standardized by the W3C [3], support only downward navigation in XML trees. As thoroughly discussed in [19], it is mainly due to the discrepancy between the XML data model and the type model, namely regular tree types [2].

Since XPath backward axes are the main source of difficulty, one may want to translate XPath selection queries with backward axes into equivalent queries with only forward axes. Olteanu et al. [39,40] propose such translations which generate a query containing the same number of joins, *i.e.*, identity-based equality, as reverse steps or a query without joins but of exponential size. The translations, however, are defined only for XPath and it is unclear how to extend them to deal with XQuery-like languages. Møller et al. [41] propose static typechecking for XSLT [42] programs which builds on a context-sensitive flow analysis. Although they introduce some approximation for abstract evaluation of XPath axes with respect to DTDs, they experimentally validate using a number of benchmarks that their algorithm is highly precise. Benzaken et al. develop a precise type inference system for XPath in their work on type-based XML projection [43]. Their system handles backward axes and is also sound and complete for a particular class of regular tree types that are *\**-guarded, non-recursive, and parent-unambiguous. In contrast, our inference system for XPath axes is exact with no such restrictions on types.

Benedikt and Cheney [44] propose a type system for the XQuery Update Facility language [45] assuming the existence of a sound typechecker for XPath axes. In the work on independence analysis of XML queries and updates [46], they use satisfiability solvers [20,47,48] to decide disjointness of selection queries, which may contain backward axes. In [47,48], weak monadic second-order logic of two successors (WS2S) [49] is used, which is one of the most expressive decidable logic when both regular tree types and XML queries are considered. However, the satisfiability problem for WS2S is known to be non-elementary. Our work is based on the tree logic and its associated satisfiability solver used in [20]. The main difference is that while [20] considers only XPath, we consider a core fragment of XQuery including element construction. Moreover, while in [20] values are defined as sets of nodes, in this work they are defined as sequences of nodes which may come from different trees and also retain their original tree context independently for navigation.

Kobayashi et al. [22] propose higher-order multi-parameter tree transducers (HMTTs) and study their verification problems, where input and output specifications are given as regular tree languages and check if the trees generated by a given HMTT conform to the output specification, given input trees satisfying the input specification. They translate an HMTT verification problem into a model checking problem for higher-order recursion schemes with finite data domains, which can be solved by using an extension of the model checking algorithm for recursion schemes proposed in [50]. Their algorithm is sound but incomplete for general HMTTs. It is, however, both sound and complete for linear HMTTs which traverse each input tree at most once. Although the complexity of the algorithm is hyper-exponential, Kobayashi et al. experimentally show that practical implementation is feasible. Since linear HMTTs subsume macro tree transducers, many XML typechecking problems can also be translated into linear HMTT verification problems (through several intermediate steps). However, it is unclear and should be further investigated how to address those programs involving XPath backward axes based on the HMTT verification method.

Recently, Castagna et al. [17] and Genevès and Gesbert [19] independently propose an extended type language to describe not only a given XML tree node but also its context. In [17], the authors extend the core calculus of CDuce [34] with zipper data structures [18], which denote the position in the surrounding tree of the value they annotate as well as its current path from the root. By annotating not only values but also types with zippers, they allow tree navigation in any direction and typecheck such navigational expressions precisely (in their work, zipped values and zipped types play a similar role as focused trees and formula-enriched sequence types, respectively). Unlike our type system, however, their typechecking is not exact for XPath axes. For example, when  $\$v$  is of type `element * {<A/>, <B/>, <C/>}`, the type of  $\$v/child::*$  is deduced as an approximate type `(<A/> | <B/> | <C/>)*` instead of an exact type `(<A/>, <B/>, <C/>)`. More precisely, in [17], a sequence type  $(\rho_1, \dots, \rho_n)$  is approximated into a more general repetition type of the form  $(\rho_1 | \dots | \rho_n)^*$ . In contrast, when  $\$v/child::*$  is of type `(<A/>, <B/>, <C/>)`, our backward type system infers an exact type `element * {<A/>, <B/>, <C/>}` for  $\$v$ . Meanwhile, Castagna et al. [17] also propose a translation from XQuery 3.0 core [7,51], with newly added value and type case analysis and higher-order functions, into the extended CDuce and provide a type system for XQuery 3.0 via the translation. In contrast to [17], currently we do not support function declarations and applications, and thus higher-order functions as well. However, because regular tree types extended with arrow types can be translated into tree logic formulas and their subtype relation can be decided through the logic's decision procedure [52], we expect that our type system can be extended with higher-order functions at least in theory.

This work builds on our previous work [19] which proposes the idea of using focused trees to denote XML values and of combining regular tree types with tree logic formulas to describe both tree nodes and their contexts simultaneously, and thus supports all the major navigational features of the XQuery core. The main difference is that while we use forward inference in [19], we use backward inference in this work. Our backward type inference is arguably more complex because we need to analyze the structure of the output type as well as the given expression (in particular, inference rules for for-loop expressions are simpler in [19]), but as a trade-off it provides an exact typechecking algorithm for XPath axes (typechecking for XPath axes is not exact in [19]). Another difference is that while we use a small-step operational semantics for the XQuery core in [19], we use a denotational semantics in this work because it is more suitable for proving properties of our backward type inference. Considering all these aspects, it would be quite interesting to combine the two approaches.

## 7. Conclusion

In this paper, we propose a novel backward type inference system for XQuery as an alternative method to forward type inference. Specifically, the contributions of the paper are summarized as follows. First, we define a focused-tree-based denotational semantics for a navigational fragment of XQuery, including all major XPath axes. Second, we propose a novel

tree-logic-based backward type inference system for XPath axes and prove its soundness and completeness. In contrast to ours, forward type inference is only sound. Finally, based on this result, we propose a sound backward type inference system for the XQuery core, with a characterized complexity.

An interesting direction for future work would be to develop a bidirectional typechecking algorithm by combining both backward and forward type inference methods. The basic idea is to typecheck for-loop expressions using forward type inference, thus obtaining a lower complexity than our backward approach, while typechecking XPath axes using backward type inference, thus obtaining better precision than the forward approaches such as in [17,19]. In doing so, one possible difficulty would be to find minimal type annotations to enable effective bidirectional typechecking.

## Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgement

The authors thank the anonymous reviewers for their helpful comments. This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIP) (No. 2016R1C1B1015095) and by the ANR project CLEAR (ANR-16-CE25-0010).

## References

- [1] S. Boag, D. Chamberlin, M. Fernández, D. Florescu, J. Robie, J. Siméon, XQuery 1.0: An XML Query Language, second edition, W3C Recommendation, December 2010, <https://www.w3.org/TR/2010/REC-xquery-20101214/>.
- [2] H. Hosoya, J. Vouillon, B.C. Pierce, Regular expression types for XML, *ACM Trans. Program. Lang. Syst.* 27 (1) (2005) 46–90, <https://doi.org/10.1145/1053468.1053470>.
- [3] D. Draper, M. Dyck, P. Fankhauser, M. Fernández, A. Malhotra, K. Rose, M. Rys, J. Siméon, P. Wadler, XQuery 1.0 and XPath 2.0 Formal Semantics, second edition, W3C Recommendation, December 2010, <http://www.w3.org/TR/xquery-semantics/>.
- [4] J. Clark, S. DeRose, XML Path Language (XPath) Version 1.0, W3C Recommendation, November 1999, <https://www.w3.org/TR/1999/REC-xpath-19991116/>.
- [5] A. Berglund, S. Boag, D. Chamberlin, M.F. Fernández, M. Kay, J. Robie, J. Siméon, XML Path Language (XPath) 2.0, second edition, W3C Recommendation, December 2010, <http://www.w3.org/TR/xpath20>.
- [6] J. Robie, D. Chamberlin, M. Dyck, J. Snelson, XML Path Language (XPath) 3.0, W3C Recommendation, April 2014, <http://www.w3.org/TR/xpath-30/>.
- [7] J. Robie, D. Chamberlin, M. Dyck, J. Snelson, XQuery 3.0: An XML Query Language, W3C Recommendation, April 2014, <http://www.w3.org/TR/xquery-30/>.
- [8] D. Colazzo, C. Sartiani, Precision and complexity of XQuery type inference, in: *Proceedings of the 13th International ACM SIGPLAN Symposium on Principles and Practices of Declarative Programming, PPDP '11, 2011*, pp. 89–100.
- [9] A. Tozawa, Towards static type checking for XSLT, in: *ACM Symposium on Document Engineering, 2001*, pp. 18–27.
- [10] T. Milo, D. Suciu, V. Vianu, Typechecking for XML transformers, *J. Comput. Syst. Sci.* 66 (1) (2003) 66–97, [https://doi.org/10.1016/S0022-0000\(02\)00030-2](https://doi.org/10.1016/S0022-0000(02)00030-2).
- [11] T. Perst, H. Seidl, Macro forest transducers, *Inf. Process. Lett.* 89 (3) (2004) 141–149, <https://doi.org/10.1016/j.ipl.2003.05.001>.
- [12] S. Maneth, A. Berlea, T. Perst, H. Seidl, XML type checking with macro tree transducers, in: *Proceedings of the Twenty-Fourth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS '05, 2005*, pp. 283–294.
- [13] A. Tozawa, XML type checking using high-level tree transducer, in: *Proceedings of the 8th International Symposium on Functional and Logic Programming, FLOPS '06, 2006*, pp. 81–96.
- [14] A. Frisch, H. Hosoya, Towards practical typechecking for macro tree transducers, in: *Proceedings of the 11th International Conference on Database Programming Languages, DBPL'07, Springer-Verlag, 2007*, pp. 246–260.
- [15] J. Engelfriet, H. Vogler, Macro tree transducers, *J. Comput. Syst. Sci.* 31 (1) (1985) 71–146, [https://doi.org/10.1016/0022-0000\(85\)90066-2](https://doi.org/10.1016/0022-0000(85)90066-2).
- [16] J. Engelfriet, The time complexity of typechecking tree-walking tree transducers, *Acta Inform.* 46 (2) (2009) 139–154, <https://doi.org/10.1007/s00236-008-0087-y>.
- [17] C. Castagna, H. Im, K. Nguyen, V. Benzaken, A core calculus for XQuery 3.0: combining navigational and pattern matching approaches, in: *Proceedings of the 24th European Symposium on Programming, 2015*, pp. 232–256.
- [18] G. Huet, The zipper, *J. Funct. Program.* 7 (5) (1997) 549–554, <https://doi.org/10.1017/S0956796897002864>.
- [19] P. Genevès, N. Gesbert, XQuery and static typing: tackling the problem of backward axes, in: *Proceedings of the 20th ACM SIGPLAN International Conference on Functional Programming, 2015*, pp. 88–100.
- [20] P. Genevès, N. Layaïda, A. Schmitt, Efficient static analysis of XML paths and types, in: *Proceedings of the 28th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '07, 2007*, pp. 342–351.
- [21] S. Maneth, T. Perst, H. Seidl, Exact XML type checking in polynomial time, in: *Proceedings of the 11th International Conference on Database Theory, ICDT'07, Springer-Verlag, Berlin, Heidelberg, 2006*, pp. 254–268.
- [22] N. Kobayashi, N. Tabuchi, H. Unno, Higher-order multi-parameter tree transducers and recursion schemes for program verification, in: *Proceedings of the 37th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '10, ACM, New York, NY, USA, 2010*, pp. 495–508.
- [23] N. Alon, T. Milo, F. Neven, D. Suciu, V. Vianu, XML with data values: typechecking revisited, *J. Comput. Syst. Sci.* 66 (4) (2003) 688–727, [https://doi.org/10.1016/S0022-0000\(03\)00032-1](https://doi.org/10.1016/S0022-0000(03)00032-1).
- [24] H. Hosoya, B.C. Pierce, XDuce: a statically typed XML processing language, *ACM Trans. Internet Technol.* 3 (2) (2003) 117–148, <https://doi.org/10.1145/767193.767195>.
- [25] M. Murata, D. Lee, M. Mani, K. Kawaguchi, Taxonomy of XML schema languages using formal language theory, *ACM Trans. Internet Technol.* 5 (4) (2005) 660–704, <https://doi.org/10.1145/1111627.1111631>.
- [26] R.M. Amadio, L. Cardelli, Subtyping recursive types, *ACM Trans. Program. Lang. Syst.* 15 (4) (1993) 575–631, <https://doi.org/10.1145/155183.155231>.
- [27] A. Brüggemann-Klein, D. Wood, One-unambiguous regular languages, *Inf. Comput.* 140 (2) (1998) 229–253, <https://doi.org/10.1006/inco.1997.2688>.

- [28] P. Genevès, N. Layaïda, A. Schmitt, N. Gesbert, Efficiently deciding  $\mu$ -calculus with converse over finite trees, *ACM Trans. Comput. Log.* 16 (2) (2015) 16, <https://doi.org/10.1145/2724712>.
- [29] M. Wand, A simple algorithm and proof for type inference, *Fundam. Inform.* 10 (1987) 115–122.
- [30] F. Pottier, Simplifying subtyping constraints: a theory, *Inf. Comput.* 170 (2) (2001) 153–183, <https://doi.org/10.1006/inco.2001.2963>.
- [31] G.M. Bierman, A.D. Gordon, C. Hrițcu, D. Langworthy, Semantic subtyping with an SMT solver, *J. Funct. Program.* 22 (1) (2012) 31–105, <https://doi.org/10.1017/S0956796812000032>.
- [32] A. Frisch, G. Castagna, V. Benzaken, Semantic subtyping: dealing set-theoretically with function, union, intersection, and negation types, *J. ACM* 55 (4) (2008) 19, <https://doi.org/10.1145/1391289.1391293>.
- [33] S. Heubach, T. Mansour, *Combinatorics of Compositions and Words*, Chapman and Hall/CRC, 2009.
- [34] V. Benzaken, G. Castagna, A. Frisch, CDuce: an XML-centric general-purpose language, in: *Proceedings of the Eighth ACM SIGPLAN International Conference on Functional Programming, ICFP '03*, 2003, pp. 51–63.
- [35] A. Möller, M.I. Schwartzbach, The design space of type checkers for XML transformation languages, in: *Proceedings of the 10th International Conference on Database Theory, ICDT '05*, Springer-Verlag, 2005, pp. 17–36.
- [36] V. Benzaken, G. Castagna, H. Hosoya, B.C. Pierce, S. Vansummeren, XML typechecking, in: *Encyclopedia of Database Systems*, Springer US, 2009, pp. 3646–3650.
- [37] S. Hakuta, S. Maneth, K. Nakano, H. Iwasaki, XQuery streaming by forest transducers, in: *IEEE 30th International Conference on Data Engineering, ICDE '14*, 2014, pp. 952–963.
- [38] A. Aho, J. Ullman, Translations on a context free grammar, *Inf. Control* 19 (5) (1971) 439–475, [https://doi.org/10.1016/S0019-9958\(71\)90706-6](https://doi.org/10.1016/S0019-9958(71)90706-6).
- [39] D. Olteanu, H. Meuss, T. Furche, F. Bry, XPath: looking forward, in: *XML-Based Data Management and Multimedia Engineering – EDBT 2002 Workshops*, 2002, pp. 109–127.
- [40] D. Olteanu, Forward node-selecting queries over trees, *ACM Trans. Database Syst.* 32 (1) (2007), <https://doi.org/10.1145/1206049.1206052>.
- [41] A. Möller, M.O. Olesen, M.I. Schwartzbach, Static validation of XSL transformations, *ACM Trans. Program. Lang. Syst.* 29 (4) (2007), <https://doi.org/10.1145/1255450.1255454>.
- [42] J. Clark, XSL Transformations (XSLT) Version 1.0, W3C Recommendation, November 1999, <https://www.w3.org/TR/1999/REC-xslt-19991116>.
- [43] V. Benzaken, G. Castagna, D. Colazzo, K. Nguyen, Optimizing XML querying using type-based document projection, *ACM Trans. Database Syst.* 38 (1) (2013) 4.
- [44] M. Benedikt, J. Cheney, Semantics, types and effects for XML updates, in: *Proceedings of the 12th International Symposium on Database Programming Languages, DBPL '09*, Springer-Verlag, 2009, pp. 1–17.
- [45] J. Robie, D. Chamberlin, M. Dyck, D. Florescu, J. Melton, J. Siméon, XQuery Update Facility 1.0, W3C Recommendation, March 2011, <http://www.w3.org/TR/xquery-update-10/>.
- [46] M. Benedikt, J. Cheney, Destabilizers and independence of XML updates, *Proc. VLDB Endow.* 3 (1–2) (2010) 906–917, <https://doi.org/10.14778/1920841.1920956>.
- [47] N. Klarlund, A. Möller, *MONA Version 1.4 User Manual*, BRICS, January 2001.
- [48] P. Genevès, N. Layaïda, Deciding XPath containment with MSO, *Data Knowl. Eng.* 63 (1) (2007) 108–136, <https://doi.org/10.1016/j.datak.2006.11.003>.
- [49] J. Doner, Tree acceptors and some of their applications, *J. Comput. Syst. Sci.* 4 (5) (1970) 406–451, [https://doi.org/10.1016/S0022-0000\(70\)80041-1](https://doi.org/10.1016/S0022-0000(70)80041-1).
- [50] N. Kobayashi, Model-checking higher-order functions, in: *Proceedings of the 11th ACM SIGPLAN Conference on Principles and Practice of Declarative Programming, PDP '09*, ACM, New York, NY, USA, 2009, pp. 25–36.
- [51] M. Benedikt, H. Vu, Higher-order functions and structured datatypes, in: *Proceedings of the 15th International Workshop on the Web and Databases, WebDB 2012*, 2012, pp. 43–48.
- [52] N. Gesbert, P. Genevès, N. Layaïda, Parametric polymorphism and semantic subtyping: the logical connection, in: *Proceedings of the 16th ACM SIGPLAN International Conference on Functional Programming, ICFP '11*, 2011, pp. 107–116.