



HAL
open science

A Core Calculus for XQuery 3.0

Giuseppe Castagna, Hyeonseung Im, Kim Nguyễn, Véronique Benzaken

► **To cite this version:**

Giuseppe Castagna, Hyeonseung Im, Kim Nguyễn, Véronique Benzaken. A Core Calculus for XQuery 3.0: Combining Navigational and Pattern Matching Approaches. ESOP '15: 24th European Symposium on Programming, 2015, London, United Kingdom. 10.1007/978-3-662-46669-8_10. hal-01104872

HAL Id: hal-01104872

<https://inria.hal.science/hal-01104872>

Submitted on 19 Jan 2015

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A Core Calculus for XQuery 3.0

Combining Navigational and Pattern Matching Approaches

Giuseppe Castagna¹, Hyeonseung Im², Kim Nguyễn³, and Véronique Benzaken³

¹ CNRS, PPS, Univ. Paris Diderot, Sorbonne Paris Cité, Paris, France

² Inria, LIG, Univ. Grenoble-Alpes, Grenoble, France

³ LRI, Université Paris-Sud, Orsay, France

Abstract. XML processing languages can be classified according to whether they extract XML data by paths or patterns. The strengths of one category correspond to the weaknesses of the other. In this work, we propose to bridge the gap between these two classes by considering two languages, one in each class: XQuery (for path-based extraction) and CDuce (for pattern-based extraction). To this end, we extend CDuce so as it can be seen as a succinct core λ -calculus that captures XQuery 3.0. The extensions we consider essentially allow CDuce to implement XPath-like navigational expressions by pattern matching and precisely type them. The elaboration of XQuery 3.0 into the extended CDuce provides a formal semantics and a sound static type system for XQuery 3.0 programs.

1 Introduction

With the establishment of XML as a standard for data representation and exchange, a wealth of XML-oriented programming languages have emerged. They can be classified into two distinct classes according to whether they extract XML data by applying paths or patterns. The strengths of one class correspond to the weaknesses of the other. In this work, we propose to bridge the gap between these classes and to do so we consider two languages each representing a distinct class: XQuery and CDuce.

XQuery [23] is a declarative language standardized by the W3C that relies heavily on XPath [21,22] as a data extraction primitive. Interestingly, the latest version of XQuery (version 3.0, very recently released [25]) adds several functional traits: type and value case analysis and functions as first-class citizens. However, while the W3C specifies a standard for document types (XML Schema [26]), it says little about the typing of XQuery programs (the XQuery 3.0 recommendation goes as far as saying that static typing is “implementation defined” and hence optional). This is a step back from the XQuery 1.0 Formal Semantics [24] which gives sound (but sometime imprecise) typing rules for XQuery.

In contrast, CDuce [4], which is used in production but issued from academic research, is a statically-typed functional language with, in particular, higher-order functions and powerful pattern matching tailored for XML data. Its key characteristic is its type algebra, which is based on *semantic subtyping* [10] and features recursive types, type constructors (product, record, and arrow types)

XQuery code

```

1  declare function get_links($page, $print) {
2      for $i in $page/descendant::a[not(ancestor::b)]
3          return $print($i)
4  }
5  declare function pretty($link) {
6      typeswitch($link)
7      case $l as element(a)
8          return switch ($l/@class)
9              case "style1"
10                 return <a href={$l/@href}><b>{$l/text()}</b></a>
11                 default return $l
12      default return $link
13  }

```

CDuce code

```

14 let get_links (page: <_>) (print: <a>_ -> <a>_) : [ <a>_ * ] =
15     match page with
16     | <a>_ & x -> [ (print x) ]
17     | < (_\`b) > l -> (transform l with (i & <_>) -> get_links i print)
18     | _ -> [ ]
19 let pretty (<a>_ -> <a>_ ; Any\<a>_ -> Any\<a>_)
20     | <a class="style1" href=h ..> l -> <a href=h>[ <b>l ]
21     | x -> x

```

Fig. 1: Document transformation in XQuery 3.0 and CDuce

and general Boolean connectives (union, intersection, and negation of types) as well as singleton types. This type algebra is particularly suited to express the types of XML documents and relies on the same foundation as the one that underpins XML Schema: regular tree languages. Moreover, the CDuce type system not only supports *ad-hoc* polymorphism (through overloading and subtyping) but also has recently been extended with parametric polymorphism [5,6].

Figure 1 highlights the key features as well as the shortcomings of both languages by defining the same two functions *get_links* and *pretty* in each language. Firstly, *get_links* (*i*) takes an XHTML document *\$page* and a function *\$print* as input, (*ii*) computes the sequence of all hypertext links (a-labelled elements) of the document that do not occur below a bold element (b-labelled elements), and (*iii*) applies the *print* argument to each link in the sequence, returning the sequence of the results. Secondly, *pretty* takes anything as argument and performs a case analysis. If the argument is a link whose `class` attribute has the value "style1", the output is a link with the same target (`href` attribute) and whose text is embedded in a bold element. Otherwise, the argument is unchanged.

We first look at the *get_links* function. In XQuery, collecting every “a” element of interest is straightforward: it is done by the XPath expression at Line 2:

$$\$page/descendant::a[not(ancestor::b)]$$

In a nutshell, an XPath expression is a sequence of steps that (i) select sets of nodes along the specified axis (here `descendant` meaning the descendants of the root node of *\$page*), (ii) keep only those nodes in the axis that have a particular label (here “a”), and (iii) further filter the results according to a Boolean condition (here `not(ancestor::b)` meaning that from a candidate “a” node, the step `ancestor::b` must return an empty result). At Lines 2–3, the `for_return` expression binds in turn each element of the result of the XPath expression to the variable *\$i*, evaluates the `return` expression, and concatenates the results. Note that there is no type annotation and that this function would fail at runtime if *\$page* is not an XML element or if *\$print* is not a function.

In clear contrast, in the CDuce program, the interface of *get_links* is fully specified (Line 14). It is curried and takes two arguments. The first one is *page* of type `<_>_`, which denotes any XML element (`_` denotes a wildcard pattern and is a synonym of the type `Any`, the type of all values, while `<s>t` is the type of an XML element with tag of type *s* and content of type *t*). The second argument is *print* of type `<a>_ → <a>_`, which is the type of functions that take an “a” element (whose content is anything) and return an “a” element. The final output is a value of type `[<a>_*]`, which denotes a possibly empty sequence of “a” elements (in CDuce’s types, the content of a sequence is described by a regular expression on types). The implementation of *get_links* in CDuce is quite different from its XQuery counterpart: following the functional idiom, it is defined as a recursive function that traverses its input recursively and performs a case analysis through pattern matching. If the input is an “a” element (Line 16), it binds the input to the capture variable *x*, evaluates *print x*, and puts the result in a sequence (denoted by square brackets). If the input is an XML element whose tag is *not* `b` (“\” stands for difference, so `_b` matches any value different from `b`)⁴, it captures the content of the element (a sequence) in *l* and applies itself recursively to each element of *l* using the `transform_with` construct whose behavior is the same as XQuery’s `for`. Lastly, if the result is not an element (or it is a “b” element), it stops the recursion and returns the empty sequence.

For the *pretty* function (which is inspired from the example given in §3.16.2 of the XQuery 3.0 recommendation [25]), the XQuery version (Lines 5–13) first performs a “type switch”, which tests whether the input *\$link* has label `a`. If so, it extracts the value of the `class` attribute using an XPath expression (Line 8) and performs a case analysis on that value. In the case where the attribute is `style1`, it re-creates an “a” element (with a nested “b” element) extracting the relevant part of the input using XPath expressions. The CDuce version (Lines 19–21) behaves in the same way but collapses all the cases in a single pattern matching. If the input is an “a” element with the desired `class` attribute, it binds the contents of the `href` attribute and the element to the variables *h* and *l*, respectively (the `..` matches possible further attributes), and builds the desired output; otherwise, the input is returned unchanged. Interestingly, this function is *overloaded*. Its signature is composed of two arrow types: if the input is an “a” element, so is the output; if the input is something else than an “a” element, so

⁴ In CDuce, one has to use `b` in conjunction with `\` to denote XML tag `b`.

is the output (& in types and patterns stands for intersection). Note that it is safe to use the *pretty* function as the second argument of the *get_links* function since $\langle a \rangle_{-} \rightarrow \langle a \rangle_{-} \& (\text{Any} \setminus \langle a \rangle_{-} \rightarrow \text{Any} \setminus \langle a \rangle_{-})$ is a subtype of $\langle a \rangle_{-} \rightarrow \langle a \rangle_{-}$ (an intersection is always smaller than or equal to the types that compose it).

Here we see that the strength of one language is the weakness of the other: CDuce provides static typing, a fine-grained type algebra, and a pattern matching construct that cleanly unifies type and value case analysis. XQuery provides through XPath a declarative way to navigate a document, which is more concise and less brittle than using hand-written recursive functions (in particular, at Line 16 in the CDuce code, there is an implicit assumption that a link cannot occur below another link; the recursion stops at “a” elements).

Contributions. The main contribution of the paper is to unify the navigational and pattern matching approaches and to define a formal semantics and type system of XQuery 3.0. Specifically, we extend CDuce so as it can be seen as a succinct core λ -calculus that can express XQuery 3.0 programs as follows.

First, we allow one to navigate in CDuce values, both downward and upward. A natural way to do so in a functional setting is to use *zippers à la Huet* [18] to annotate values. Zippers denote the position in the surrounding tree of the value they annotate as well as its current path from the root. We extend CDuce not only with zipped values (*i.e.*, values annotated by zippers) but also with *zipped types*. By doing so, we show that we can navigate not only in any direction in a document but also in a *precisely typed* way, allowing one to express constraints on the path in which a value is within a document.

Second, we extend CDuce pattern matching with accumulating variables that allow us to encode *recursive* XPath axes (such as **descendant** and **ancestor**). It is well known that typing such recursive axes goes well beyond regular tree languages and that approximations in the type system are needed. Rather than giving ad-hoc built-in functions for **descendant** and **ancestor**, we define the notion of *type operators* and parameterize the CDuce type system (and dynamic semantics) with these operators. Soundness properties can then be shown in a modular way without hard-coding any specific typing rules in the language. With this addition, XPath navigation can be encoded simply in CDuce’s pattern matching constructs and it is just a matter of syntactic sugar definition to endow CDuce with nice declarative navigational expressions such as those successfully used in XQuery or XSLT.

The last (but not least) step of our work is to define a “normal form” for XQuery 3.0 programs, extending both the original XQuery Core normal form of [24] and its recent adaptation to XQuery 3.0 (dubbed XQ_H) proposed by Benedikt and Vu [3]. In this normal form, navigational (*i.e.*, structural) expressions are well separated from data value expressions (ordering, node identity testing, *etc.*). We then provide a translation from XQuery 3.0 Core to CDuce extended with navigational patterns. The encoding provides for free an effective and efficient typechecking algorithm for XQuery 3.0 programs (described in Figure 9 of Section 5.1) as well as a formal and compact specification of their semantics. Even more interestingly, it provides a solid formal basis to start further studies on the

Pre-values	$w ::= c \mid (w, w) \mid \mu f^{(t \rightarrow t; \dots; t \rightarrow t)}(x).e$
Zippers	$\delta ::= \bullet \mid \mathbf{L}(w)_\delta \cdot \delta \mid \mathbf{R}(w)_\delta \cdot \delta$
Values	$v ::= w \mid (v, v) \mid (w)_\delta$
Expressions	$e ::= v \mid x \mid \dot{x} \mid (e, e) \mid (e)_\bullet \mid o(e, \dots, e)$ $\mid \text{match } e \text{ with } p \rightarrow e \mid p \rightarrow e$
Pre-types	$u ::= b \mid c \mid u \times u \mid u \rightarrow u \mid u \vee u \mid \neg u \mid 0$
Zipper types	$\tau ::= \bullet \mid \top \mid \mathbf{L}(u)_\tau \cdot \tau \mid \mathbf{R}(u)_\tau \cdot \tau \mid \tau \vee \tau \mid \neg \tau$
Types	$t ::= u \mid t \times t \mid t \rightarrow t \mid t \vee t \mid \neg t \mid (u)_\tau$
Pre-patterns	$q ::= t \mid x \mid \dot{x} \mid (q, q) \mid q \mid q \mid q \& q \mid (x := c)$
Zipper patterns	$\varphi ::= \tau \mid \mathbf{L}p \cdot \varphi \mid \mathbf{R}p \cdot \varphi \mid \varphi \mid \varphi$
Patterns	$p ::= q \mid (p, p) \mid p \mid p \mid p \& p \mid (q)_\varphi$

Fig. 2: Syntax of expressions, types, and patterns

definition of XQuery 3.0 and its properties. *A minima*, it is straightforward to use this basis to add overloaded functions to XQuery (*e.g.*, to give a precise type to *pretty*). More crucially, the recent advances on polymorphism for semantic subtyping [5,6,7] can be transposed to this basis to provide a polymorphic type system and type inference algorithm both to XQuery 3.0 and to the extended CDuce language defined here. Polymorphic types are the missing ingredient to make higher-order functions yield their full potential and to remove any residual justification of the absence of standardization of the XQuery 3.0 type system.

Plan. Section 2 presents the core typed λ -calculus equipped with zipper-annotated values, accumulators, constructors, recursive functions, and pattern matching. Section 3 gives its semantics, type system, and the expected soundness property. Section 4 turns this core calculus into a full-fledged language using several syntactic constructs and encodings. Section 5 uses this language as a compilation target for XQuery. Lastly, Section 6 compares our work to other related approaches and concludes. Proofs and some technical definitions are given in Appendix.

2 Syntax

We extend the CDuce language [4] with zippers *à la* Huet [18]. To ensure the well-foundedness of the definition, we stratify it, introducing first pre-values (which are standard CDuce values) and then values, which are pre-values possibly indexed by a zipper; we proceed similarly for types and patterns. The definition is summarized in Figure 2. Henceforth we denote by \mathcal{V} the set of all values and by Ω a special value that represents runtime error and does not inhabit any type. We also denote by \mathcal{E} and \mathcal{T} the set of all expressions and all types, respectively.

2.1 Values and Expressions

Pre-values (ranged over by w) are the usual CDuce values without zipper annotations. Constants are ranged over by c and represent integers (1, 2, ...),

characters ('a', 'b', ...), atoms ('nil', 'true', 'false', 'foo', ...), etc. A value (w, w) represents pairs of pre-values. Our calculus also features recursive functions (hence the μ binder instead of the traditional λ) with explicit, overloaded types (the set of types that index the recursion variable, forming the *interface* of the function). Values (ranged over by v) are pre-values, pairs of values, or pre-values annotated with a *zipper* (ranged over by δ). Zippers are used to record the path covered when traversing a data structure. Since the product is the only construct, we need only three kinds of zippers: the empty one (denoted by \bullet) which intuitively denotes the starting point of our navigation, and two zippers $L(w)_\delta \cdot \delta$ and $R(w)_\delta \cdot \delta$ which denote respectively the path to the left and right projection of a pre-value w , which is itself reachable through δ . To ease the writing of several zipper related functions, we chose to record in the zipper the whole “stack” of values we have visited (each tagged with a left or right indication), instead of just keeping the unused component as is usual.

Example 1. Let v be the value $((1, (2, 3)))_\bullet$. Its first projection is the value $(1)_{L((1, (2, 3)))_\bullet} \cdot \bullet$ and its second projection is the value $((2, 3))_{R((1, (2, 3)))_\bullet} \cdot \bullet$, the first projection of which being $(2)_{L((2, 3))_{R((1, (2, 3)))_\bullet} \cdot \bullet} \cdot R((1, (2, 3)))_\bullet \cdot \bullet$.

As one can see in this example, keeping values in the zipper (instead of pre-values) seems redundant since the same value occurs several times (see how δ is duplicated in the definition of zippers). The reason for this duplication is purely syntactic: it makes the writing of types and patterns that match such values much shorter (intuitively, to go “up” in a zipper, it is only necessary to extract the previous value while keeping it un-annotated — *i.e.*, having $Lw \cdot \delta$ in the definition instead of $L(w)_\delta \cdot \delta$ — would require a more complex treatment to reconstruct the parent). We also stress that zipped values are meant to be used only for internal representation: the programmer will be allowed to write just pre-values (not values or expressions with zippers) and be able to obtain and manipulate zippers only by applying CDuce functions and pattern matching (as defined in the rest of the paper) and never directly.

Expressions include values (as previously defined), variables (ranged over by x, y, \dots), accumulators (which are a particular kind of variables, ranged over by \dot{x}, \dot{y}, \dots), and pairs. An expression $(e)_\bullet$ annotates e with the empty zipper \bullet . The pattern matching expression is standard (with a first match policy) and will be thoroughly presented in Section 3. Our calculus is parameterized by a set \mathcal{O} of built-in operators ranged over by o . Before describing the use of operators and the set of operators defined in our calculus (in particular the operators for projection and function application), we introduce our type algebra.

2.2 Types

We first recall the CDuce type algebra, as defined in [10], where types are interpreted as sets of values and the subtyping relation is semantically defined by using this interpretation (*i.e.*, $\llbracket t \rrbracket = \{v \mid \vdash v : t\}$ and $s \leq t \stackrel{\text{def}}{\iff} \llbracket s \rrbracket \subseteq \llbracket t \rrbracket$).

Pre-types u (as defined in Figure 2) are the usual CDuce types, which are possibly infinite terms with two additional requirements:

1. (regularity) the number of distinct subterms of u is finite;
2. (contractiveness) every infinite branch of u contains an infinite number of occurrences of either product types or function types.

We use b to range over basic types (`int`, `bool`, ...). A singleton type c denotes the type that contains only the constant value c . The empty type \emptyset contains no value. Product and function types are standard: $u_1 \times u_2$ contains all the pairs (w_1, w_2) for $w_i \in u_i$, while $u_1 \rightarrow u_2$ contains all the (pre-)value functions that when applied to a value in u_1 , if such application terminates then it returns a value in u_2 . We also include type connectives for union and negation (intersections are encoded below) with their usual set-theoretic interpretation. Infiniteness of pre-types accounts for recursive types and regularity implies that pre-types are finitely representable, for instance, by recursive equations or by the explicit μ -notation. Contractiveness [2] excludes both ill-formed (*i.e.*, unguarded) recursions such as $\mu X.X$ as well as meaningless type definitions such as $\mu X.X \vee X$ or $\mu X.\neg X$ (unions and negations are finite). Finally, subtyping is defined as set-theoretic containment (u_1 is a subtype of u_2 , denoted by $u_1 \leq u_2$, if all values in u_1 are also in u_2) and it is decidable in EXPTIME (see [10]).

A *zipper type* τ is a possibly infinite term that is regular as for pre-types and contractive in the sense that every infinite branch of τ must contain an infinite number of occurrences of either left or right projection. The singleton type \bullet is the type of the empty zipper and \top denotes the type of all zippers, while $L(u)_\tau \cdot \tau$ (resp., $R(u)_\tau \cdot \tau$) denotes the type of zippers that encode the left (resp., right) projection of some value of pre-type u . We use $\tau_1 \wedge \tau_2$ to denote $\neg(\neg\tau_1 \vee \neg\tau_2)$.

The type algebra of our core calculus is then defined as pre-types possibly indexed by zipper types. As for pre-types, a *type* t is a possibly infinite term that is both regular and contractive. We write $t \wedge s$ for $\neg(\neg t \vee \neg s)$, $t \setminus s$ for $t \wedge \neg s$, and $\mathbb{1}$ for $\neg\emptyset$; in particular, $\mathbb{1}$ denotes the super-type of all types (it contains all values). We also define the following notations (we use \equiv both for syntactic equivalence and definition of syntactic sugar):

- $\mathbb{1}_{\text{prod}} \equiv \mathbb{1} \times \mathbb{1}$ the super-type of all product types
- $\mathbb{1}_{\text{fun}} \equiv \emptyset \rightarrow \mathbb{1}$ the super-type of all arrow types
- $\mathbb{1}_{\text{basic}} \equiv \mathbb{1} \setminus (\mathbb{1}_{\text{prod}} \vee \mathbb{1}_{\text{fun}} \vee (\mathbb{1})_\top)$ the super-type of all basic types
- $\mathbb{1}_{\text{NZ}} \equiv \mu X.(X \times X) \vee (\mathbb{1}_{\text{basic}} \vee \mathbb{1}_{\text{fun}})$ the type of all pre-values (*i.e.*, Not Zipped)

It is straightforward to extend the subtyping relation of pre-types (*i.e.*, the one defined in [10]) to our types: the addition of $(u)_\tau$ corresponds to the addition of a new type constructor (akin to \rightarrow and \times) to the type algebra. Therefore, it suffices to define the interpretation of the new constructor to complete the definition of the subtyping relation (defined as containment of the interpretations). In particular, $(u)_\tau$ is interpreted as the set of all values $(w)_\delta$ such that $\vdash w : u$ and $\vdash \delta : \tau$ (both typing judgments are defined in Appendix B.1). From this we deduce that $(\mathbb{1})_\top$ (equivalently, $(\mathbb{1}_{\text{NZ}})_\top$) is the type of all (pre-)values decorated with a zipper. The formal definition is more involved (see Appendix A) but the intuition is simple: a type $(u_1)_{\tau_1}$ is a subtype of $(u_2)_{\tau_2}$ if $u_1 \leq u_2$ and τ_2 is a prefix (modulo type equivalence and subtyping) of τ_1 . The prefix containment

translates the intuition that the more we know about the context surrounding a value, the more numerous are the situations in which it can be safely used. For instance, in XML terms, if we have a function that expects an element whose parent’s first child is an integer, then we can safely apply this function to an element whose type indicates that its parent’s first child has type (a subtype of) integer *and* that its grandparent is, say, tagged by **a**.

Finally, as for pre-types, the subtyping relation for types is decidable in EXPTIME. This is easily shown by producing a straightforward linear encoding of zipper types and zipper values in pre-types and pre-values, respectively (the encoding is given in Definition 16 in Appendix A).

2.3 Operators and Accumulators

As previously explained, our calculus includes accumulators and is parameterized by a set \mathcal{O} of operators. These have the following formal definitions:

Definition 2 (Operator). *An operator is a 4-tuple $(o, n_o, \overset{\circ}{\rightsquigarrow}, \overset{\circ}{\rightarrow})$ where o is the name (symbol) of the operator, n_o is its arity, $\overset{\circ}{\rightsquigarrow} \subseteq \mathcal{V}^{n_o} \times \mathcal{E} \cup \{\Omega\}$ is its reduction relation, and $\overset{\circ}{\rightarrow} : \mathcal{T}^{n_o} \rightarrow \mathcal{T}$ is its typing function.*

In other words, an operator is an applicative symbol, equipped with both a dynamic (\rightsquigarrow) and a static (\rightarrow) semantics. The reason for making $\overset{\circ}{\rightsquigarrow}$ a relation is to account for non-deterministic operators (*e.g.*, random choice). Note that an operator may fail, thus returning the special value Ω during evaluation.

Definition 3 (Accumulator). *An accumulator \hat{x} is a variable equipped with a binary operator $Op(\hat{x}) \in \mathcal{O}$ and initial value $Init(\hat{x}) \in \mathcal{V}$.*

2.4 Patterns

Now that we have defined types and operators, we can define patterns. Intuitively, patterns are types with capture variables that are used either to extract subtrees from an input value or to test its “shape”. As before, we first recall the definition of standard CDuce patterns (here called pre-patterns), enrich them with accumulators, and then extend the whole with zippers.

A pre-pattern q , as defined in Figure 2, is either a type constraint t , or a capture variable x , or an accumulator \hat{x} , or a pair (q_1, q_2) , or an alternative $q_1 \mid q_2$, or a conjunction $q_1 \ \& \ q_2$, or a default case $(x := c)$. It is a possibly infinite term that is regular as for pre-types and contractive in the sense that every infinite branch of q must contain an infinite number of occurrences of pair patterns. Moreover, the subpatterns forming conjunctions must have distinct capture variables and those forming alternatives the same capture variables. A *zipper pattern* φ is a possibly infinite term that is both regular and contractive as for zipper types. Finally, a pattern p is a possibly infinite term with the same requirements as pre-patterns. Besides, the subpatterns q and φ forming a zipper pattern $(q)_\varphi$ must have distinct capture variables. We denote by $\text{Var}(p)$ the set of capture variables occurring in p and by $\text{Acc}(p)$ the set of accumulators occurring in p .

$$\begin{array}{c}
E ::= [] \mid (E, e) \mid (e, E) \mid (E)_{\bullet} \mid \text{match } E \text{ with } p_1 \rightarrow e_1 \mid p_2 \rightarrow e_2 \mid o(e, \dots, E, \dots, e) \\
\\
\frac{(v_1, \dots, v_{n_o}) \overset{o}{\rightsquigarrow} e}{o(v_1, \dots, v_{n_o}) \rightsquigarrow e} \quad \frac{\{\dot{x} \mapsto \text{Init}(\dot{x}) \mid \dot{x} \in \text{Acc}(p_1)\}; \square \vdash v/p_1 \rightsquigarrow \sigma, \gamma}{\text{match } v \text{ with } p_1 \rightarrow e_1 \mid p_2 \rightarrow e_2 \rightsquigarrow e_1[\sigma; \gamma]} \\
\frac{\{\dot{x} \mapsto \text{Init}(\dot{x}) \mid \dot{x} \in \text{Acc}(p_1)\}; \square \vdash v/p_1 \rightsquigarrow \Omega \quad \{\dot{x} \mapsto \text{Init}(\dot{x}) \mid \dot{x} \in \text{Acc}(p_2)\}; \square \vdash v/p_2 \rightsquigarrow \sigma, \gamma}{\text{match } v \text{ with } p_1 \rightarrow e_1 \mid p_2 \rightarrow e_2 \rightsquigarrow e_2[\sigma; \gamma]} \\
\\
\frac{e \rightsquigarrow e'}{E[e] \rightsquigarrow E[e']} \quad \frac{}{e \rightsquigarrow \Omega} \quad \left(\begin{array}{l} \text{if no other rule applies} \\ \text{and } e \text{ is not a value} \end{array} \right)
\end{array}$$

Fig. 3: Operational semantics (reduction contexts and rules)

3 Semantics

In this section, the most technical one, we present the operational semantics and the type system of our calculus, and state the expected soundness properties.

3.1 Operational Semantics

We define a call-by-value, small-step operational semantics for our core calculus, using the reduction contexts and reduction rules given in Figure 3, where Ω is a special value representing a runtime error.

Of course, most of the actual semantics is hidden (the careful reader will have noticed that applications and projections are not explicitly included in the syntax of our expressions). Most of the work happens either in the semantics of operators or in the matching v/p of a value v against a pattern p . Such a matching, if it succeeds (*i.e.*, if it does not return Ω), returns two substitutions, one (ranged over by γ) from the capture variables of p to values and the other (ranged over by δ) from the accumulators to values. These two substitutions are simultaneously applied (noted $e_i[\sigma; \gamma]$) to the expression e_i of the pattern p_i that succeeds, according to a first match policy (v/p_2 is evaluated only if v/p_1 fails). Before explaining how to derive the pattern matching judgments “ $_ \vdash v/p \rightsquigarrow _$ ” (in particular, the meaning of the context on the LHS of the turnstile “ \vdash ”), we introduce a minimal set of operators: application, projections, zipper erasure, and sequence building (we use **sans-serif** font for concrete operators). We only give their reduction relation and defer their typing relation to Section 3.2.

Function application: the operator **app**($_$, $_$) implements the usual β -reduction:

$$v, v' \overset{\text{app}}{\rightsquigarrow} e[v/f; v'/x] \quad \text{if } v = \mu f^{(\dots)}(x).e$$

and $v, v' \overset{\text{app}}{\rightsquigarrow} \Omega$ if v is not a function. As customary, $e[v/x]$ denotes the capture-avoiding substitution of v for x in e , and we write $e_1 e_2$ for **app**(e_1, e_2).

Projection: the operator $\pi_1(_)$ (resp., $\pi_2(_)$) implements the usual first (resp., second) projection for pairs:

$$(v_1, v_2) \overset{\pi_i}{\rightsquigarrow} v_i \quad \text{for } i \in \{1, 2\}$$

The application of the above operators returns Ω if the input is not a pair.

Zipper erasure: given a zipper-annotated value, it is sometimes necessary to remove the zipper (*e.g.*, to embed this value into a new data structure). This is achieved by the following remove $\text{rm}(_)$ and deep remove $\text{drm}(_)$ operators:

$$\begin{array}{ccc} (w)_\delta & \xrightarrow{\text{rm}} & w \\ v & \xrightarrow{\text{rm}} & v \quad \text{if } v \neq (w)_\delta \end{array} \qquad \begin{array}{ccc} w & \xrightarrow{\text{drm}} & w \\ (w)_\delta & \xrightarrow{\text{drm}} & w \\ (v_1, v_2) & \xrightarrow{\text{drm}} & (\text{drm}(v_1), \text{drm}(v_2)) \end{array}$$

The former operator only erases the top-level zipper (if any), while the latter erases all zippers occurring in its input.

Sequence building: given a sequence (encoded *à la* Lisp) and an element, we define the operators $\text{cons}(_)$ and $\text{snoc}(_)$ that insert an input value at the beginning and at the end of the input sequence:

$$\begin{array}{ccc} v, v' & \xrightarrow{\text{cons}} & (v, v') \\ v, \text{'nil} & \xrightarrow{\text{snoc}} & (v, \text{'nil}) \\ v, (v', v'') & \xrightarrow{\text{snoc}} & (v', \text{snoc}(v, v'')) \end{array}$$

The applications of these operators yield Ω on other inputs.

To complete our presentation of the operational semantics, it remains to describe the semantics of pattern matching. Intuitively, when matching a value v against a pattern p , subparts of p are recursively applied to corresponding subparts of v until a base case is reached (which is always the case since all values are finite). As usual, when a pattern variable is confronted with a subvalue, the binding is stored as a substitution. We supplement this usual behavior of pattern matching with two novel features. First, we add *accumulators*, that is, special variables in which results are accumulated during the recursive matching. The reason for keeping these two kinds of variables distinct is explained in Section 3.2 and is related to type inference for patterns. Second, we parameterize pattern matching by a zipper of the current value so that it can properly update the zipper when navigating the value (which should be of the pair form).

These novelties are reflected by the semantics of pattern matching, which is given by the judgment $\sigma; \delta^? \vdash v/p \rightsquigarrow \sigma', \gamma$, where v is a value, p a pattern, γ a mapping from $\text{Var}(p)$ to values, and σ and σ' are mappings from accumulators to values. $\delta^?$ is an optional zipper value, which is either δ or a none value \square (we consider $(v)_\square$ to be v). The judgment “returns” the result of matching the value v against the pattern p (noted v/p), that is, two substitutions: γ for capture variables and σ' for accumulators. Since the semantics is given compositionally, the matching may happen on a subpart of an “outer” matched value. Therefore, the judgment records on the LHS of the turnstile the context of the outer value explored so far: σ stores the values already accumulated during the matching, while $\delta^?$ tracks the possible zipper of the outer value (or it is \square if the outer value has no zipper). The context is “initialized” in the two rules of the operational semantics of `match` in Figure 3, by setting each accumulator of the pattern to its initial value (function `Init()`) and the outer zipper to \square .

Judgments for pattern matching are derived by the rules given in Figure 4. The rules `pat-acc`, `pat-pair-zip`, and `zpat-*` are novel, as they extend pattern matching with accumulators and zippers, while the others are derived from [4,9].

$$\begin{array}{c}
\frac{(\vdash v : t)}{\sigma; \delta^? \vdash v/t \rightsquigarrow \sigma, \emptyset} \text{ pat-type} \quad \frac{}{\sigma; \delta^? \vdash v/\dot{x} \rightsquigarrow \sigma[\text{Op}(\dot{x})(v_{\delta^?}, \sigma(\dot{x}))/\dot{x}], \emptyset} \text{ pat-acc} \\
\frac{}{\sigma; \delta^? \vdash v/x \rightsquigarrow \sigma, \{x \mapsto v_{\delta^?}\}} \text{ pat-var} \quad \frac{}{\sigma; \delta^? \vdash v/(x := c) \rightsquigarrow \sigma, \{x \mapsto c\}} \text{ pat-def} \\
\frac{\sigma; \square \vdash v_1/p_1 \rightsquigarrow \sigma', \gamma_1 \quad \sigma'; \square \vdash v_2/p_2 \rightsquigarrow \sigma'', \gamma_2}{\sigma; \square \vdash (v_1, v_2)/(p_1, p_2) \rightsquigarrow \sigma'', \gamma_1 \oplus \gamma_2} \text{ pat-pair} \\
\frac{\sigma; \mathbf{L}(w_1, w_2)_\delta \cdot \delta \vdash w_1/p_1 \rightsquigarrow \sigma', \gamma_1 \quad \sigma'; \mathbf{R}(w_1, w_2)_\delta \cdot \delta \vdash w_2/p_2 \rightsquigarrow \sigma'', \gamma_2}{\sigma; \delta \vdash (w_1, w_2)/(p_1, p_2) \rightsquigarrow \sigma'', \gamma_1 \oplus \gamma_2} \text{ pat-pair-zip} \\
\frac{\sigma; \delta^? \vdash v/p_1 \rightsquigarrow \sigma', \gamma}{\sigma; \delta^? \vdash v/p_1 \mid p_2 \rightsquigarrow \sigma', \gamma} \text{ pat-or1} \quad \frac{\sigma; \delta^? \vdash v/p_1 \rightsquigarrow \Omega \quad \sigma; \delta^? \vdash v/p_2 \rightsquigarrow \sigma', \gamma}{\sigma; \delta^? \vdash v/p_1 \mid p_2 \rightsquigarrow \sigma', \gamma} \text{ pat-or2} \\
\frac{\sigma; \delta^? \vdash v/p_1 \rightsquigarrow \sigma', \gamma_1 \quad \sigma'; \delta^? \vdash v/p_2 \rightsquigarrow \sigma'', \gamma_2}{\sigma; \delta^? \vdash v/p_1 \& p_2 \rightsquigarrow \sigma'', \gamma_1 \oplus \gamma_2} \text{ pat-and} \\
\frac{\sigma; \delta \vdash w/q \rightsquigarrow \sigma', \gamma_1 \quad \sigma' \vdash \delta/\varphi \rightsquigarrow \sigma'', \gamma_2}{\sigma; \square \vdash (w)_\delta/(q)_\varphi \rightsquigarrow \sigma'', \gamma_1 \oplus \gamma_2} \text{ pat-zip} \quad \frac{(\vdash \delta : \tau)}{\sigma \vdash \delta/\tau \rightsquigarrow \sigma, \emptyset} \text{ zpat-type} \\
\frac{\sigma; \square \vdash (w)_\delta/p \rightsquigarrow \sigma', \gamma_1 \quad \sigma' \vdash \delta/\varphi \rightsquigarrow \sigma'', \gamma_2 \quad \gamma = \gamma_1 \oplus \gamma_2}{\sigma \vdash \mathbf{L}(w)_\delta \cdot \delta/\mathbf{L}p \cdot \varphi \rightsquigarrow \sigma'', \gamma} \text{ zpat-left} \quad \frac{\sigma; \square \vdash (w)_\delta/p \rightsquigarrow \sigma', \gamma_1 \quad \sigma' \vdash \delta/\varphi \rightsquigarrow \sigma'', \gamma_2 \quad \gamma = \gamma_1 \oplus \gamma_2}{\sigma \vdash \mathbf{R}(w)_\delta \cdot \delta/\mathbf{R}p \cdot \varphi \rightsquigarrow \sigma'', \gamma} \text{ zpat-right} \\
\frac{\sigma \vdash \delta/\varphi_1 \rightsquigarrow \sigma', \gamma}{\sigma \vdash \delta/\varphi_1 \mid \varphi_2 \rightsquigarrow \sigma', \gamma} \text{ zpat-or1} \quad \frac{\sigma \vdash \delta/\varphi_1 \rightsquigarrow \Omega \quad \sigma \vdash \delta/\varphi_2 \rightsquigarrow \sigma', \gamma}{\sigma \vdash \delta/\varphi_1 \mid \varphi_2 \rightsquigarrow \sigma', \gamma} \text{ zpat-or2} \\
\frac{(\text{otherwise})}{\sigma; \delta^? \vdash v/p \rightsquigarrow \Omega} \text{ pat-error} \quad \frac{(\text{otherwise})}{\sigma \vdash \delta/\varphi \rightsquigarrow \Omega} \text{ zpat-error}
\end{array}$$

where $\gamma_1 \oplus \gamma_2 \stackrel{\text{def}}{=} \begin{array}{l} \{x \mapsto \gamma_1(x) \mid x \in \text{dom}(\gamma_1) \setminus \text{dom}(\gamma_2)\} \\ \cup \{x \mapsto \gamma_2(x) \mid x \in \text{dom}(\gamma_2) \setminus \text{dom}(\gamma_1)\} \\ \cup \{x \mapsto (\gamma_1(x), \gamma_2(x)) \mid x \in \text{dom}(\gamma_1) \cap \text{dom}(\gamma_2)\} \end{array}$

Fig. 4: Pattern matching

There are three base cases for matching: testing the input value against a type (rule **pat-type**), updating the environment σ for accumulators (rule **pat-acc**), or producing a substitution γ for capture variables (rules **pat-var** and **pat-def**). Matching a pattern (p_1, p_2) only succeeds if the input is a pair and the matching of each subpattern against the corresponding subvalue succeeds (rule **pat-pair**). Furthermore, if the value being matched was below a zipper (*i.e.*, the current zipper context is a δ and not—as in **pat-pair**— \square), we update the current zipper context (rule **pat-pair-zip**); notice that in this case the matched value must be a pair of pre-values since zipped values cannot be nested. An alternative pattern $p_1 \mid p_2$ first tries to match the pattern p_1 and if it fails, tries the pattern p_2 (rules **pat-or1** and **pat-or2**). The matching of a conjunction pattern $p_1 \& p_2$ succeeds if and only if the matching of both patterns succeeds (rule **pat-and**). For a zipper constraint $(q)_\varphi$, the matching succeeds if and only if the input value is annotated by a zipper, *e.g.*, $(w)_\delta$, and both the matching of w with q and δ with φ succeed

(rule **pat-zip**). It requires the zipper context to be \square since we do not allow nested zipped values. When matching w with q , we record the zipper δ into the context so that it can be updated (in the rule **pat-pair-zip**) while navigating the value.

The matching of a zipper pattern φ against a zipper δ (judgments $\sigma \vdash \delta/\varphi \rightsquigarrow \sigma', \gamma$ derived by the **zpat-*** rules) is straightforward: it succeeds if both φ and δ are built using the same constructor (either **L** or **R**) and the componentwise matching succeeds (rules **zpat-left** and **zpat-right**). If the zipper pattern is a zipper type, the matching tests the input zipper against the zipper type (rule **zpat-type**), and alternative zipper patterns $\varphi_1 | \varphi_2$ follow the same first match policy as alternative patterns. If none of the rules is applicable, the matching fails (rules **pat-error** and **zpat-error**). Note that initially the environment σ contains $\text{Init}(\dot{x})$ for each accumulator \dot{x} in $\text{Acc}(p)$ (rules for **match** in Figure 3).

Intuitively, γ is built when returning from the recursive descent in p , while σ is built using a *fold*-like computation. It is the typing of such fold-like computations that justifies the addition of accumulators (instead of relying on plain functions). But before presenting the type system of the language, we illustrate the behavior of pattern matching by some examples.

Example 4. Let $v \equiv (2, ('true, (3, 'nil)))$, $\text{Init}(\dot{x}) = 'nil$, $\text{Op}(\dot{x}) = \text{cons}$, and $\sigma \equiv \{\dot{x} \mapsto 'nil\}$. Then, we have the following matchings:

1. $\sigma; \square \vdash v/(\text{int}, (x, _)) \rightsquigarrow \emptyset, \{x \mapsto 'true\}$
2. $\sigma; \square \vdash v/\mu X.((x \& \text{int} | _, X) | (x := 'nil)) \rightsquigarrow \emptyset, \{x \mapsto (2, (3, 'nil))\}$
3. $\sigma; \square \vdash v/\mu X.((\dot{x}, X) | 'nil) \rightsquigarrow \{\dot{x} \mapsto (3, ('true, (2, 'nil)))\}, \emptyset$

In the first case, the input v (the sequence [2 'true 3] encoded *à la* Lisp) is matched against a pattern that checks if the first element has type **int** (rule **pat-type**), binds the second element to x (rule **pat-var**), and ignores the rest of the list (rule **pat-type**, since the anonymous variable “ $_$ ” is just an alias for $\mathbb{1}$).

The second case is more involved since the pattern is recursively defined. Because of the first match policy of rule **pat-or1**, the product part of the pattern is matched recursively until the atom **'nil** is reached. When that is the case, the variable x is bound to a default value **'nil**. When returning from this recursive matching, since x occurs both on the left and on the right of the product (in $x \& \text{int}$ and in X itself), a pair of the binding found in each part is formed (third set in the definition of \oplus in Figure 4), thus yielding a mapping $\{x \mapsto (3, 'nil)\}$. Returning again from the recursive call, only the “ $_$ ” part of the pattern matches the input **'true** (since it is not of type **int**, the intersection test fails). Therefore, the binding for this step is only the binding for the right part (second case of the definition of \oplus). Lastly, when reaching the top-level pair, $x \& \text{int}$ matches 2 and a pair is formed from this binding and the one found in the recursive call, yielding the final binding $\{x \mapsto (2, (3, 'nil))\}$.

The third case is more intuitive. The pattern just recurses the input value, calling the accumulation function for \dot{x} along the way for each value against which it is confronted. Since the operator associated with \dot{x} is **cons** (which builds a pair of its two arguments) and the initial value is **'nil**, this has the effect of computing the reversal of the list.

Note the key difference between the second and third case. In both cases, the structure of the pattern (and the input) dictates the traversal, but in the second case, it also dictates *how* the binding is built (if v was a tree and not a list, the binding for x would also be a tree in the second case). In the third case, the way the binding is built is defined by the semantics of the operator and independent of the input. This allows us to reverse sequences or flatten tree structures, both of which are operations that escape the expressiveness of regular tree languages/regular patterns, but which are both necessary to encode XPath.

3.2 Type System

The main difficulty in devising the type system is to type pattern matching and, more specifically, to infer the types of the accumulators occurring in patterns.

Definition 5 (Accepted input of an operator). *The accepted input of an operator $(o, n, \overset{\circ}{\sim}, \overset{\circ}{\rightarrow})$ is the set $\mathbb{I}(o)$, defined as:*

$$\mathbb{I}(o) = \{(v_1, \dots, v_n) \in \mathcal{V}^n \mid ((v_1, \dots, v_n) \overset{\circ}{\sim} e) \wedge (e \rightsquigarrow^* v) \Rightarrow v \neq \Omega\}$$

Definition 6 (Exact input). *An operator o has an exact input if and only if $\mathbb{I}(o)$ is (the interpretation of) a type.*

We can now state a first soundness theorem, which characterizes the set of all values that make a given pattern succeed:

Theorem 7 (Accepted types). *Let p be a pattern such that for every \hat{x} in $\text{Acc}(p)$, $\text{Op}(\hat{x})$ has an exact input. Then, the set of all values v such that $\{\hat{x} \mapsto \text{Init}(\hat{x}) \mid \hat{x} \in \text{Acc}(p)\}; \square \vdash v/p \not\rightsquigarrow \Omega$ is a type. We call this set the accepted type of p and denote it by $\downarrow p$.*

We next define the type system for our core calculus, in the form of a judgment $\Gamma \vdash e : t$ which states that in a typing environment Γ (i.e., a mapping from variables and accumulators to types) an expression e has type t . This judgment is derived by the set of rules given in Figure 10 in Appendix. Here, we show only the most important rules, namely those for accumulators and zippers:

$$\frac{}{\Gamma \vdash \hat{x} : \Gamma(\hat{x})} \quad \frac{\vdash w : t \quad \vdash \delta : \tau \quad t \leq \mathbb{1}_{\text{NZ}}}{\Gamma \vdash (w)_\delta : (t)_\tau} \quad \frac{\vdash e : t \quad t \leq \mathbb{1}_{\text{NZ}}}{\Gamma \vdash (e)_\bullet : (t)_\bullet}$$

which rely on an auxiliary judgment $\vdash \delta : \tau$ stating that a zipper δ has zipper type τ . The rule for operators is:

$$\frac{\forall i = 1..n_o, \Gamma \vdash e_i : t_i \quad t_1, \dots, t_{n_o} \overset{\circ}{\rightarrow} t}{\Gamma \vdash o(e_1, \dots, e_{n_o}) : t} \text{ for } o \in \mathcal{O}$$

which types operators using their associated typing function. Last but not least, the rule to type pattern matching expressions is:

$$\frac{\begin{array}{l} t \leq \downarrow p_1 \vee \downarrow p_2 \\ t_1 \equiv t \wedge \downarrow p_1 \quad t_2 \equiv t \wedge \neg \downarrow p_1 \\ \Sigma_i \equiv \{\hat{x} \mapsto \text{Init}(\hat{x}) \mid \hat{x} \in \text{Acc}(p_i)\} \end{array} \quad \begin{array}{l} \Gamma \vdash e : t \\ \Gamma_i \equiv \square \vdash t_i/p_i \quad \Gamma'_i \equiv \Sigma_i; \square \vdash t_i//p_i \\ \Gamma \cup \Gamma_i \cup \Gamma'_i \vdash e_i : t'_i \end{array}}{\Gamma \vdash \text{match } e \text{ with } p_1 \rightarrow e_1 \mid p_2 \rightarrow e_2 : \bigvee_{\{i \mid t_i \neq \emptyset\}} t'_i} \quad (i = 1, 2)$$

This rule requires that the type t of the matched expression is smaller than $\wr p_1 \wr \vee \wr p_2 \wr$ (*i.e.*, the set of all values accepted by any of the two patterns), that is, that the matching is exhaustive. Then, it accounts for the first match policy by checking e_1 in an environment inferred from values produced by e and that match p_1 ($t_1 \equiv t \wedge \wr p_1 \wr$) and by checking e_2 in an environment inferred from values produced by e and that *do not* match p_1 ($t_2 \equiv t \wedge \neg \wr p_1 \wr$). If one of these branches is unused (*i.e.*, if $t_i \simeq \emptyset$ where \simeq denotes semantic equivalence, that is, $\leq \cap \geq$), then its type does not contribute to the type of the whole expression (*cf.* §4.1 of [4] to see why, in general, this must not yield an “unused case” error). Each right-hand side e_i is typed in an environment enriched with the types for capture variables (computed by $\square \sim t_i / p_i$) and the types for accumulators (computed by $\Sigma_i; \square \sim t_i // p_i$). While the latter is specific to our calculus, the former is standard except it is parameterized by a zipper type as for the semantics of pattern matching (its precise computation is described in [9] and already implemented in the CDuce compiler except the zipper-related part: see Figure 11 in Appendix for the details). As before, we write $\tau^?$ to denote an optional zipper type, *i.e.*, either τ or a none type \square , and consider $(t)_{\square}$ to be t .

To compute the types of the accumulators of a pattern p when matched against a type t , we first initialize an environment Σ by associating each accumulator \dot{x} occurring in p with the singleton type for its initial value $\text{Init}(\dot{x})$ ($\Sigma_i \equiv \{\dot{x} \mapsto \text{Init}(\dot{x}) \mid \dot{x} \in \text{Acc}(p_i)\}$). The type environment is then computed by generating a set of mutually recursive equations where the important ones are (see Figure 12 in Appendix for the complete definition):

$$\begin{aligned}
\Sigma; \tau^? \sim t // \dot{x} &= \Sigma[s/\dot{x}] && \text{if } (t)_{\tau^?}, \Sigma(\dot{x}) \xrightarrow{\text{Op}(\dot{x})} s \\
\Sigma; \tau^? \sim t // p_1 \mid p_2 &= \Sigma; \tau^? \sim t // p_1 && \text{if } t \leq \wr p_1 \wr \\
\Sigma; \tau^? \sim t // p_1 \mid p_2 &= \Sigma; \tau^? \sim t // p_2 && \text{if } t \leq \neg \wr p_1 \wr \\
\Sigma; \tau^? \sim t // p_1 \mid p_2 &= (\Sigma; \tau^? \sim (t \wedge \wr p_1 \wr) // p_1) \sqcup (\Sigma_1; \tau^? \sim (t \wedge \neg \wr p_1 \wr) // p_2) && \text{otherwise}
\end{aligned}$$

When an accumulator \dot{x} is matched against a type t , the type of the accumulator is updated in Σ , by applying the typing function of the operator associated with \dot{x} to the type $(t)_{\tau^?}$ and the type computed thus far for \dot{x} , namely $\Sigma(\dot{x})$. The other equations recursively apply the matching on the subcomponents while updating the zipper type argument $\tau^?$ and merge the results using the “ \sqcup ” operation. This operation implements the fact that if an accumulator \dot{x} has type t_1 in a subpart of a pattern p and type t_2 in another subpart (*i.e.*, both subparts match), then the type of \dot{x} is the union $t_1 \vee t_2$.

The equations for computing the type environment for accumulators might be *not* well-founded. Both patterns and types are possibly infinite (regular) terms and therefore one has to guarantee that the set of generated equations is finite. This depends on the typing of the operators used for the accumulators. Before stating the termination condition (as well as the soundness properties of the type system), we give the typing functions for the operators we defined earlier.

Function application: it is typed by computing the minimum type satisfying the following subtyping relation: $s, t \xrightarrow{\text{app}} \min\{t' \mid s \leq t \rightarrow t'\}$, provided that $s \leq t \rightarrow \mathbb{1}$ (this min always exists and is computable: see [10]).

Projection: to type the first and second projections, we use the property that if $t \leq \mathbb{1} \times \mathbb{1}$, then t can be decomposed in a finite union of product types (we use Π_i to denote the set of the i -th projections of these types: see Lemma 19 in Appendix B for the formal definition): $t \xrightarrow{\pi_i} \bigvee_{s \in \Pi_i(t)} s$, provided that $t \leq \mathbb{1} \times \mathbb{1}$.

Zipper erasure: the top-level erasure $\xrightarrow{\text{rm}}$ simply removes the top-level zipper type annotation, while the deep erasure $\xrightarrow{\text{drm}}$ is typed by recursively removing the zipper annotations from the input type. Their precise definition can be found in Appendix B.4.

Sequence building: it is typed in the following way:

$$\begin{array}{l} t_1, \text{'nil} \xrightarrow{\text{cons}} \mu X.((t_1 \times X) \vee \text{'nil}) \\ t_1, \mu X.((t_2 \times X) \vee \text{'nil}) \xrightarrow{\text{cons}} \mu X.(((t_1 \vee t_2) \times X) \vee \text{'nil}) \\ \\ t_1, \text{'nil} \xrightarrow{\text{snoc}} \mu X.((t_1 \times X) \vee \text{'nil}) \\ t_1, \mu X.((t_2 \times X) \vee \text{'nil}) \xrightarrow{\text{snoc}} \mu X.(((t_1 \vee t_2) \times X) \vee \text{'nil}) \end{array}$$

Notice that the output types are approximations: the operator “cons()” is *less* precise than returning a pair of two values since, for instance, it approximates any sequence type by an infinite one (meaning that any information on the length of the sequence is lost) and approximates the type of all the elements by a single type which is the union of all the elements (meaning that the information on the order of elements is lost). As we show next, this loss of precision is instrumental in typing accumulators and therefore pattern matching.

Example 8. Consider the matching of a pattern p against a value v of type t :

$$\begin{array}{l} p \equiv \mu X.((\hat{x} \& (\text{'a} \mid \text{'b})) \mid \text{'nil} \mid (X, X)) \\ v \equiv (\text{'a}, ((\text{'a}, (\text{'nil}, (\text{'b}, \text{'nil}))), (\text{'b}, \text{'nil}))) \\ t \equiv \mu Y.((\text{'a} \times (Y \times (\text{'b} \times \text{'nil}))) \vee \text{'nil}) \end{array}$$

where $\text{Op}(\hat{x}) = \text{snoc}$ and $\text{Init}(\hat{x}) = \text{'nil}$. We have the following matching and type environment:

$$\begin{array}{l} \{\hat{x} \mapsto \text{'nil}\}; \square \vdash v/p \rightsquigarrow \{\hat{x} \mapsto (\text{'a}, (\text{'a}, (\text{'b}, (\text{'b}, \text{'nil}))))\}, \emptyset \\ \{\hat{x} \mapsto \text{'nil}\}; \square \vdash t/p = \{\hat{x} \mapsto \mu Z.(((\text{'a} \vee \text{'b}) \times Z) \vee \text{'nil})\} \end{array}$$

Intuitively, with the usual sequence notation (precisely defined in Section 4), v is nothing but the nested sequence [‘a [‘a [‘a [‘b] ‘b] and pattern matching just flattens the input sequence, binding \hat{x} to [‘a ‘a ‘b ‘b]. The type environment for \hat{x} is computed by recursively matching each product type in t with the pattern (X, X) , the singleton type ‘a or ‘b with $\hat{x} \& (\text{'a} \mid \text{'b})$, and ‘nil with ‘nil. Since the operator associated with \hat{x} is `snoc` and the initial type is ‘nil, when \hat{x} is matched against ‘a for the first time, its type is updated to $\mu Z.((\text{'a} \times Z) \vee \text{'nil})$. Then, when \hat{x} is matched against ‘b, its type is updated

to the final output type which is the encoding of $[('a \vee 'b)*]$. Here, the approximation in the typing function for `snoc` is important because the exact type of \hat{x} is the union for $n \in \mathbb{N}$ of $['a^n 'b^n]$, that is, the sequences of `'a`'s followed by the same number of `'b`'s, which is beyond the expressivity of regular tree languages.

We conclude this section with statements for type soundness of our calculus (see Appendix C for more details).

Definition 9 (Sound operator). *An operator $(o, n, \overset{o}{\rightsquigarrow}, \overset{o}{\rightarrow})$ is sound if and only if $\forall v_1, \dots, v_{n_o} \in \mathcal{V}$ such that $\vdash v_1 : t_1, \dots, \vdash v_{n_o} : t_{n_o}$, if $t_1, \dots, t_{n_o} \overset{o}{\rightarrow} s$ and $v_1, \dots, v_{n_o} \overset{o}{\rightsquigarrow} e$ then $\vdash e : s$.*

Theorem 10 (Type preservation). *If all operators in the language are sound, then typing is preserved by reduction, that is, if $e \rightsquigarrow e'$ and $\vdash e : t$, then $\vdash e' : t$. In particular, $e' \neq \Omega$.*

Theorem 11. *The operators `app`, π_1 , π_2 , `drm`, `rm`, `cons`, and `snoc` are sound.*

4 Surface Language

In this section, we define the “surface” language, which extends our core calculus with several constructs:

- Sequence expressions, regular expression types and patterns
- Sequence concatenation and iteration
- XML types, XML document fragment expressions
- XPath-like patterns

While most of these traits are syntactic sugar or straightforward extensions, we took special care in their design so that: (i) they cover various aspects of XML programming and (ii) they are expressive enough to encode a large fragment of XQuery 3.0.

Sequences: we first add sequences to expressions

$$e ::= \dots \mid [e \cdots e]$$

where a sequence expression denotes its encoding *à la* Lisp, that is, $[e_1 \cdots e_n]$ is syntactic sugar for $(e_1, (\dots, (e_n, \text{nil})))$.

Regular expression types and patterns: regular expressions over types and patterns are defined as

$$\begin{aligned} \text{(Regexp. over types)} \quad R & ::= t \mid R \mid R \mid RR \mid R* \mid \epsilon \\ \text{(Regexp. over patterns)} \quad r & ::= p \mid r \mid r \mid rr \mid r* \mid \epsilon \end{aligned}$$

with the usual syntactic sugar: $R? \equiv R \mid \epsilon$ and $R+ \equiv RR*$ (likewise for regexps on patterns). We then extend the grammar of types and patterns as follows:

$$t ::= \dots \mid [R] \quad p ::= \dots \mid [r]$$

Regular expression types are encoded using recursive types (similarly for regular expression patterns). For instance, $[\text{int}* \text{bool}?]$ can be rewritten into the recursive type $\mu X. \text{'nil} \vee (\text{bool} \times \text{'nil}) \vee (\text{int} \times X)$.

Sequence concatenation is added to the language in the form of a binary infix operator $_ @ _$ defined by:

$$\begin{array}{l} \langle \mathbf{nil}, v \rangle \xrightarrow{\textcircled{e}} v \\ \langle v_1, v_2 \rangle, v \xrightarrow{\textcircled{e}} (v_1, v_2 @ v) \end{array} \quad [R_1], [R_2] \xrightarrow{\textcircled{e}} [R_1 R_2]$$

Note that this operator is sound but cannot be used to accumulate in patterns (since it does not guarantee the termination of type environment computation). However, it has an exact typing.

Sequence iteration is added to iterate transformations over sequences without resorting to recursive functions. This is done by a family of “transform”-like operators $\text{trs}_{p_1, p_2, e_1, e_2}(_)$, indexed by the patterns and expressions that form the branches of the transformation (we omit trs ’s indexes in $\xrightarrow{\text{trs}}$):

$$\begin{array}{l} \langle \mathbf{nil} \rangle \xrightarrow{\text{trs}} \langle \mathbf{nil} \rangle \\ \langle v_1, v_2 \rangle \xrightarrow{\text{trs}} (\text{match } v_1 \text{ with } p_1 \rightarrow e_1 \mid p_2 \rightarrow e_2) @ \text{trs}_{p_1, p_2, e_1, e_2}(v_2) \end{array}$$

Intuitively, the construct “transform e with $p_1 \rightarrow e_1 \mid p_2 \rightarrow e_2$ ” iterates all the “branches” over each element of the sequence e . Each branch may return a sequence of results which is concatenated to the final result (in particular, a branch may return “ \mathbf{nil} ” to delete elements that match a particular pattern).

XML types, patterns, and document fragments: XML types (and thus patterns) can be represented as a pair of the type of the label and a sequence type representing the sequence of children, annotated by the zipper that denotes the position of document fragment of that type. We denote by $\langle t_1 \rangle t_2 \tau$ the type $(t_1 \times t_2) \tau$, where $t_1 \leq \mathbb{1}_{\text{basic}}$, $t_2 \leq [\mathbb{1}^*]$, and τ is a zipper type. We simply write $\langle t_1 \rangle t_2$ when $\tau = \top$, that is, when we do not have (or do not require) any information on the zipper type. The invariant that XML values are always given with respect to a zipper must be maintained at the level of expressions. This is ensured by extending the syntax of expressions with the construct:

$$e ::= \dots \mid \langle e \rangle e$$

where $\langle e_1 \rangle e_2$ is syntactic sugar for $(e_1, \text{drm}(e_2))_{\bullet}$. The reason for this encoding is best understood with the following example:

Example 12. Consider the code:

```

1  match v with
2    ( <a>[ _ x _* ] )⊤ -> <b>[ x ]
3  | _ -> <c>[ ]
```

According to our definition of pattern matching, x is bound to the second XML child of v and retains its zipper (in the right-hand side, we could navigate from x up to v or even above if v is not the root). However, when x is embedded into another document fragment, the zipper must be erased so that accessing the element associated with x in the *new* value can create an appropriate zipper (with respect to its new root $\langle b \rangle [\dots]$).

$$\begin{aligned}
\text{self}_0\{x \mid t\} &\equiv x \ \& \ t \mid _ \\
\text{self}\{x \mid t\} &\equiv (\text{self}_0\{x \mid t\})_{\top} \\
\text{child}\{x \mid t\} &\equiv (\langle _ \rangle [(\text{self}_0\{x \mid t\})^*] \mid _)_{\top} \\
\text{desc-or-self}_0\{x \mid t\} &\equiv \mu X. (\text{self}_0\{x \mid t\} \ \& \ \langle _ \rangle [X^*]) \mid _ \\
\text{desc-or-self}\{x \mid t\} &\equiv (\text{desc-or-self}_0\{x \mid t\})_{\top} \\
\text{desc}\{x \mid t\} &\equiv (\langle _ \rangle [(\text{desc-or-self}_0\{x \mid t\})^*] \mid _)_{\top} \\
\text{foll-sibling}\{x \mid t\} &\equiv (_)_{\text{L}} (_)_{\text{L}} [(\text{self}_0\{x \mid t\})^*]_{\top} \cdot \top \\
\text{parent}\{y \mid t\} &\equiv (_)_{\text{L}} _ \cdot \mu X. ((R \ (y \ \& \ t \mid _)_{\top} \cdot (\text{L} _ \cdot \top \mid \bullet)) \mid R _ \cdot X) \mid _ \\
\text{prec-sibling}\{y \mid t\} &\equiv (_)_{\text{L}} _ \cdot \mu X. (R \ (y \ \& \ t, _)_{\top} \cdot X) \mid (R _ \cdot (\text{L} _ \cdot \top \mid \bullet)) \mid _ \\
\text{anc}\{y \mid t\} &\equiv (_)_{\text{L}} _ \cdot \mu X. \mu Y. ((R \ (y \ \& \ t \mid _)_{\top} \cdot (\text{L} _ \cdot X \mid \bullet)) \mid R _ \cdot Y) \mid _ \\
\text{anc-or-self}\{y \mid t\} &\equiv (\text{self}\{y \mid t\} \ \& \ \text{anc}\{y \mid t\}) \mid _
\end{aligned}$$

where $\text{Op}(x) = \text{snoc}$, $\text{Init}(x) = \text{'nil}$, $\text{Op}(y) = \text{cons}$, and $\text{Init}(y) = \text{'nil}$

Fig. 5: Encoding of axis patterns

XPath-like patterns are one of the main motivations for this work. The syntax of patterns is extended as follows:

$$\begin{aligned}
(\text{Patterns}) \quad p &::= \dots \mid \text{axis}\{x \mid t\} \\
(\text{Axes}) \quad \text{axis} &::= \text{self} \mid \text{child} \mid \text{desc} \mid \text{desc-or-self} \mid \text{foll-sibling} \\
&\quad \mid \text{parent} \mid \text{anc} \mid \text{anc-or-self} \mid \text{prec-sibling}
\end{aligned}$$

The semantics of $\text{axis}\{x \mid t\}$ is to capture in x all fragments of the matched document along the axis that have type t . We show in Appendix D how the remaining two axes (**following** and **preceding**) as well as “multi-step” XPath expressions can be compiled into this simpler form. We encode axis patterns directly using recursive patterns and accumulators, as described in Figure 5. First, we remark that each pattern has a default branch “ $\dots \mid _$ ” which implements the fact that even if a pattern fails, the value is still accepted, but the default value ‘nil’ of the accumulator is returned. The so-called “downward” axes —**self**, **child**, **desc-or-self**, and **desc**— are straightforward. For **self**, the encoding checks that the matched value has type t using the auxiliary pattern self_0 , and that the value is annotated with a zipper using the zipper type annotation $(_)_{\top}$. The **child** axis is encoded by iterating self_0 on every child element of the matched value. The recursive axis **desc-or-self** is encoded using the auxiliary pattern desc-or-self_0 which matches the root of the current element (using self_0) and is recursively applied to each element of the sequence. Note the double recursion: vertically in the tree using a recursive binder and horizontally at a given level using a star. The non-reflexive variant **desc** evaluates desc-or-self_0 on every child element of the input.

The other axes heavily rely on the binary encoding of XML values and are better explained on an example. Consider the XML document and its binary tree representation given in Figure 6. The following siblings of a node (*e.g.*, $\langle c \rangle$) are reachable by inspecting the first element of the zipper, which is necessarily an **L** one. This parent is the pair representing the sequence whose tail is the sequence of following siblings (R_3 and R_2 in the figure). Applying the $\text{self}\{x \mid t\}$ axis on each element of the tail therefore filters the following siblings that are

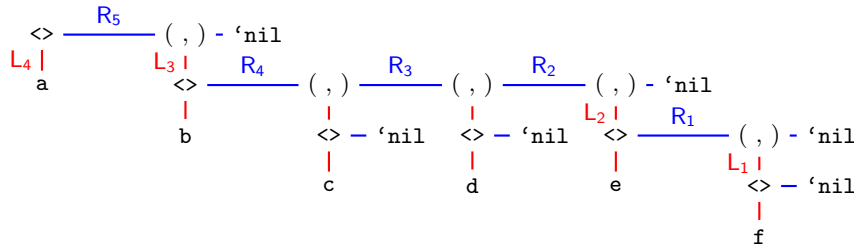


Fig. 6: A binary tree representation of an XML document
`doc = <a>[[<c>[] <d>[] <e>[<f>[]]]]`

sought (`<d>` and `<e>` in the figure). The **parent** axis is more involved. Consider for instance node `<e>`. Its parent in the XML tree can be found in the zipper associated with `<e>`. It is the last **R** component of the zipper before the next **L** component (in the figure, the zipper of `<e>` starts with L_2 , then contains its previous siblings reachable by R_2 and R_3 , and lastly its parent reachable by R_4 (which points to node ``). The encoding of the **parent** axis reproduces this walk using a recursive zipper pattern, whose base case is the last **R** before the next **L**, or the last **R** before the root (which has the empty zipper \bullet). The **prec-sibling** axis uses a similar method and collects every node reachable by **R**s and stops before the parent node (again, for node `<e>`, the preceding siblings are reached by R_2 and R_3). The **anc** axis simply iterates the parent axis recursively until there is no **L** zipper anymore (*i.e.*, until the root of the document has been reached). In the example, starting from node `<f>`, the zippers that denote the ancestors are the ones starting with an **R**, just before L_2 , L_3 , and L_4 which is the root of the document. Lastly, **anc-or-self** is simply a combination of **anc** and **self**.

For space reasons, the encoding of XPath into the navigational patterns is given in Appendix D. We just stress that, with that encoding, the CDuce version of the “*get_links*” function of the introduction becomes as compact as in XQuery:

```
let get_links (page: <_>_) (print: <a>_ -> <a>_) : [ <a>_ * ] =
  transform page/desc::a[not(anc::b)] with x -> [ (print x) ]
```

As a final remark, one may notice that patterns of forward axes use `snoc` (*i.e.*, they build the sequence of the results in order), while reverse axes use `cons` (thus reversing the results). The reason for this difference is to implement the semantics of XPath axis steps which return elements *in document order*.

5 XQuery 3.0

This section shows that our surface language can be used as a compilation target for XQuery 3.0 programs. We proceed in two steps. First, we extend the XQuery 1.0 Core fragment and XQ_H defined by Benedikt and Vu [3] to our own XQuery 3.0 Core, which we call XQ_H^+ . As with its 1.0 counterpart, XQ_H^+

1. can express *all* navigational XQuery programs, and

<code>query ::= () c <l>query</l></code>		<code>query, query x x/axis::test</code>
<code>for x in query return query</code>		<code>some x in query satisfies query</code>
<code>query(query, ..., query)</code>		<code>fun x₁ : t₁, ..., x_n : t_n as t. query</code>
<code>switch query</code>		<code>typeswitch query</code>
<code> case c return query</code>		<code> case t as x return query</code>
<code> default return query</code>		<code> default return query</code>

`test ::= node() | text() | l` (node test)

where t ranges over types and l ranges over element names.

Fig. 7: Syntax of XQ_H^+

2. explicitly separates navigational aspects from data value ones.

We later use the above separation in the translation to straightforwardly map navigational XPath expressions into extended CDuce pattern matching, and to encode data value operations (for which there can be no precise typing) by built-in CDuce functions.

5.1 XQuery 3.0 Core

Figure 7 shows the definition of XQ_H^+ , an *extension* of XQ_H . To the best of our knowledge, XQ_H was the first work to propose a “Core” fragment of XQuery which abstracts away most of the idiosyncrasies of the actual specification while retaining essential features (*e.g.*, path navigation). XQ_H^+ differs from XQ_H by the last three productions (in the yellow/gray box): it extends XQ_H with type and value cases (described informally in the introduction) and with *type annotations* on functions (which are only optional in the standard). It is well known (*e.g.*, see [24]) that full XPath expressions can be encoded using the XQuery fragment in Figure 7 (see Appendix E for an illustration).

Our translation of XQuery 3.0, defined in Figure 8, thus focuses on XQ_H^+ and has following characteristics. If one considers the “typed” version of the standard, that is, XQuery programs where function declarations have an explicit signature, then the translation to our surface language (*i*) provides a formal semantics and a typechecking algorithm for XQuery and (*ii*) enjoys the soundness property that the original XQuery programs do not yield runtime errors. In the present work, we assume that the type algebra of XQuery is the one of CDuce, rather than XMLSchema. Both share regular expression types for which subtyping is implemented as the inclusion of languages, but XMLSchema also features *nominal subtyping*. The extension of CDuce types with nominal subtyping is beyond the scope of this work and is left as future work.

In XQuery, all values are sequences: the constant “42” is considered as the *singleton sequence* that contains the element “42”. As a consequence, there are only “flat” sequences in XQuery and the only way to create nested data structures is to use XML constructs. The difficulty for our translation is thus twofold: (*i*) it needs to embed/extract values explicitly into/from sequences and (*ii*) it

$$\begin{aligned}
\llbracket () \rrbracket_{\text{XC}} &= \text{'nil} \\
\llbracket c \rrbracket_{\text{XC}} &= [c] \\
\llbracket \langle l \rangle q \langle /l \rangle \rrbracket_{\text{XC}} &= \langle [q] \rangle_{\text{XC}} \\
\llbracket q_1, q_2 \rrbracket_{\text{XC}} &= \llbracket q_1 \rrbracket_{\text{XC}} @ \llbracket q_2 \rrbracket_{\text{XC}} \\
\llbracket \$x \rrbracket_{\text{XC}} &= x \\
\llbracket \text{switch } q_1 \text{ case } c \text{ return } q_2 \text{ default return } q_3 \rrbracket_{\text{XC}} &= \text{match } \llbracket q_1 \rrbracket_{\text{XC}} \text{ with} \\
&\quad | c \rightarrow \llbracket q_2 \rrbracket_{\text{XC}} \\
&\quad | _ \rightarrow \llbracket q_3 \rrbracket_{\text{XC}} \\
\llbracket \text{typeswitch } q_1 \text{ case } t \text{ as } \$x \text{ return } q_2 \text{ default return } q_3 \rrbracket_{\text{XC}} &= \text{match } \llbracket q_1 \rrbracket_{\text{XC}} \text{ with} \\
&\quad x \& \text{seq}(t) \rightarrow \llbracket q_2 \rrbracket_{\text{XC}} \\
&\quad | _ \rightarrow \llbracket q_3 \rrbracket_{\text{XC}} \\
\llbracket \$x / \text{axis} :: \text{test} \rrbracket_{\text{XC}} &= \text{transform } x \text{ with } \text{axis}\{y \mid t(\text{test})\} \rightarrow y \\
\llbracket \text{for } \$x \text{ in } q_1 \text{ return } q_2 \rrbracket_{\text{XC}} &= \text{transform } \llbracket q_1 \rrbracket_{\text{XC}} \text{ with } x \rightarrow \llbracket q_2 \rrbracket_{\text{XC}} \\
\llbracket \text{some } \$x \text{ in } q_1 \text{ satisfies } q_2 \rrbracket_{\text{XC}} &= \text{match (transform } \llbracket q_1 \rrbracket_{\text{XC}} \text{ with} \\
&\quad x \rightarrow \text{match } \llbracket q_2 \rrbracket_{\text{XC}} \text{ with} \\
&\quad \quad [\text{'true}] \rightarrow [\text{'dummy}] \\
&\quad \quad | [\text{'false}] \rightarrow []) \\
&\quad \text{with 'nil} \rightarrow [\text{'false}] \mid _ \rightarrow [\text{'true}] \\
\llbracket \text{fun } \$x_1 : t_1, \dots, \$x_n : t_n \text{ as } t. q \rrbracket_{\text{XC}} &= \mu_{_}^{\text{seq}(t_1) \times \dots \times \text{seq}(t_n) \rightarrow \text{seq}(t)}(x_0). \\
&\quad \text{match } x_0 \text{ with } (x_1, (\dots, x_n)) \rightarrow \llbracket q \rrbracket_{\text{XC}} \\
\llbracket q(q_1, \dots, q_n) \rrbracket_{\text{XC}} &= \llbracket q \rrbracket_{\text{XC}} (\llbracket q_1 \rrbracket_{\text{XC}}, (\dots, \llbracket q_n \rrbracket_{\text{XC}}))
\end{aligned}$$

where $\text{seq}(t) \equiv (t \wedge [\mathbb{1}^*]) \vee (t \setminus [\mathbb{1}^*])$
and $t(\text{node}()) \equiv \mathbb{1}$, $t(\text{text}()) \equiv \text{String}$, $t(l) \equiv \langle l \rangle \mathbb{1}$

Fig. 8: Translation of XQ_{H}^+ into CDuce

also needs to disambiguate types: an XQuery function that takes an integer as argument can also be applied to a sequence containing only one integer.

The translation is defined by a function $\llbracket _ \rrbracket_{\text{XC}}$ that converts an XQuery query into a CDuce expression. It is straightforward and ensures that the result of a translation $\llbracket q \rrbracket_{\text{XC}}$ always has a sequence type. We assume that both languages have the same set of variables and constants. An empty sequence is translated into the atom `'nil`, a constant is translated into a singleton sequence containing that constant, and similarly for XML fragments. The sequence operator is translated into concatenation. Variables do not require any special treatment. An XPath navigation step is translated into the corresponding navigational pattern, whereas “for in” loops are encoded similarly using the `transform` construct (in XQuery, an XPath query applied to a sequence of elements is the concatenation of the individual applications). The “switch” construct is directly translated into a “match with” construct. The “typeswitch” construct works in a similar way but special care must be taken with respect to the type t that is tested. Indeed, if t is a sequence type, then its translation returns the sequence type, but if t is something else (say `int`), then it must be embedded into a sequence type. Interestingly, this test can be encoded as the CDuce type $\text{seq}(t)$ which keeps the part of t that is a sequence unchanged while embedding the part of t that is not a sequence (namely $t \setminus [\mathbb{1}^*]$) into a sequence type (*i.e.*, $[t \setminus [\mathbb{1}^*]]$). The “some $\$x$ in q_1 satisfies q_2 ” expression iterates over the sequence that is the result of the translation of q_1 , binding variable x in turn to each element, and evaluates (the translation of) q_2 in this context. If the evaluation of q_2 yields the single-

$$\begin{array}{c}
\frac{\Gamma \vdash_{\text{XQ}} q : s \quad s \leq t \quad \Gamma \vdash_{\text{XQ}} q_1 : [s*] \quad \Gamma, x : [s] \vdash_{\text{XQ}} q_2 : t \quad t \leq [\mathbb{1}*]}{\Gamma \vdash_{\text{XQ}} q : t} \quad \Gamma \vdash_{\text{XQ}} \text{for } \$x \text{ in } q_1 \text{ return } q_2 : t \\
\frac{\{y \mapsto \text{nil}\}; \square \vdash s // \text{axis}\{y \mid \mathbf{t}(\text{test})\} = \{y \mapsto t\} \quad \Gamma \vdash_{\text{XQ}} x : [s*] \quad t \leq [\mathbb{1}*] \quad t' = \min\{t' \mid t \leq [t'*]\}}{\Gamma \vdash_{\text{XQ}} \$x/\text{axis}::\text{test} : [t'*]} \text{typ-path} \\
\frac{\Gamma \vdash_{\text{XQ}} q : t \quad \begin{cases} t \not\leq \neg[c] \Rightarrow \Gamma \vdash_{\text{XQ}} q_1 : s & t_1 = s \wedge \text{seq}(t) \quad \Gamma, x : t_1 \vdash_{\text{XQ}} q_1 : t'_1 \\ t \not\leq [c] \Rightarrow \Gamma \vdash_{\text{XQ}} q_2 : s & \Gamma \vdash_{\text{XQ}} q : s \quad t_2 = s \wedge \neg\text{seq}(t) \quad \Gamma \vdash_{\text{XQ}} q_2 : t'_2 \end{cases}}{\Gamma \vdash_{\text{XQ}} \text{switch } q \quad \text{typeswitch } q} \\
\frac{\Gamma \vdash_{\text{XQ}} \text{case } c \text{ return } q_1 : s \quad \text{default return } q_2}{\Gamma \vdash_{\text{XQ}} \text{case } t \text{ as } \$x \text{ return } q_1 : \bigvee_{\{i \mid t_i \neq 0\}} t'_i} \\
\frac{\Gamma \vdash_{\text{XQ}} q_1 : [s*] \quad \Gamma, x : [s] \vdash_{\text{XQ}} q_2 : [\text{bool}]}{\Gamma \vdash_{\text{XQ}} \text{some } \$x \text{ in } q_1 \text{ satisfies } q_2 : [\text{bool}]} \\
\frac{\Gamma, x_1 : \text{seq}(t_1), \dots, x_n : \text{seq}(t_n) \vdash_{\text{XQ}} q : \text{seq}(t)}{\Gamma \vdash_{\text{XQ}} \text{fun } \$x_1 : t_1, \dots, \$x_n : t_n \text{ as } t. q : \text{seq}(t_1) \times \dots \times \text{seq}(t_n) \rightarrow \text{seq}(t)} \\
\frac{\Gamma \vdash_{\text{XQ}} q : t_1 \times \dots \times t_n \rightarrow t \quad \Gamma \vdash_{\text{XQ}} q_i : t_i \quad (i = 1..n)}{\Gamma \vdash_{\text{XQ}} q(q, \dots, q) : t}
\end{array}$$

Fig. 9: Typing rules for XQ_H^+

ton sequence **true**, then we return a dummy non-empty sequence; otherwise, we return the empty sequence. If the whole transform yields an empty sequence, it means that none of the iterated elements matched satisfied the predicate q_2 and therefore the whole expression evaluates to the singleton **false**, otherwise it evaluates to the singleton **true**. Abstractions are translated into CDuce functions, and the same treatment of “sequencing” the type is applied to the types of the arguments and type of the result. Lastly, application is translated by building nested pairs with the arguments before applying the function.

Not only does this translation ensure soundness of the original XQuery 3.0 programs, it also turns CDuce into a sandbox where one can experiment various typing features that can be readily back-ported to XQuery afterwards.

5.2 Toward and Beyond XQuery 3.0

We now discuss the salient features and address some shortcomings of XQ_H^+ . First and foremost, we can define a precise and sound type system directly on XQ_H^+ as shown in Figure 9 (standard typing rules are omitted and for the complete definition, see Appendix E). While most constructs are typed straightforwardly (the typing rules are deduced from the translation of XQ_H^+ into CDuce) it is interesting to see that the rules match those defined in XQuery Static Semantics specification [24] (with the already mentioned difference that we use CDuce types instead of XMLSchema). Two aspects however diverge from the standard. Our use of CDuce’s semantic subtyping (rather than XMLSchema’s nominal subtyping), and the rule **typ-path** where we use the formal developments of Section 3 to provide a precise typing rule for XPath navigation. Deriving the typing rules from our translation allows us to state the following theorem:

Theorem 13. *If $\Gamma \vdash_{\text{XQ}} \text{query} : t$, then $\Gamma \vdash \llbracket \text{query} \rrbracket_{\text{XC}} : t$.*

A corollary of this theorem is the soundness of the XQ_H^+ type system (since the translation of a well-typed XQ_H^+ program yields a well-typed CDuce program with the same type).

While the XQ_H^+ fragment we present here is already very expressive, it does not account for all features of XQuery. For instance, it does not feature data value comparison or sorting (*i.e.*, the `order by` construct of XQuery) nor does it account for built-in functions such as `position()`, node identifiers, and so on. However, it is known that features such as data value comparison make typechecking undecidable (see for instance [1]). We argue that the main point of this fragment is to cleanly separate structural path navigation from other data value tests for which we can add built-in operators and functions, with an hardcoded, *ad-hoc* typing rule.

Lastly, one may argue that, in practice, XQuery database engines do *not* rely on XQuery Core for evaluation but rather focus on evaluating efficiently large (multi-step, multi-predicate) XPath expressions in one go and, therefore, that normalizing XQuery programs into XQ_H^+ programs and then translating the latter into CDuce programs may seem overly naive. We show in Appendix D that XPath expressions that are purely navigational can be rewritten in a single pattern of the form: `axis{x | t}` which can then be evaluated very efficiently (that is, without performing the unneeded extra traversals of the document that a single step approach would incur).

6 Related Work and Conclusion

Our work tackles several aspects of XML programming, the salient being: (i) encoding of XPath or XPath-like expressions (including reverse axes) into regular types and patterns, (ii) recursive tree transformation using accumulators and their typing, and (iii) type systems and typechecking algorithms for XQuery.

Regarding XPath and pattern matching, the work closest to ours is the implementation of paths as patterns in XTatic. XTatic [11] is an object-oriented language featuring XDuce regular expression types and patterns [16,17]. In [12], Gapeyev and Pierce alter XDuce’s pattern matching semantics and encode a fragment of XPath as patterns. The main difference with our work is that they use a hard-coded all-match semantics (a variable can be bound to several sub-terms) to encode the accumulations of recursive axes, which are restricted by their data model to the “child” and “descendant” axes. Another attempt to use path navigation in a functional language can be found in [19] where XPath-like combinators are added to Haskell. Again, only child or descendant-like navigation is supported and typing is done in the setting of Haskell which cannot readily be applied to XML typing (results are returned as *homogeneous* sequences).

Our use of accumulators is reminiscent of Macro Tree Transducers (MTTs, [8]), that is, tree transducers (tree automata producing an output) that can also accumulate part of the input and copy it in the output. It is well known that given an input regular tree language, the type of the accumulators and results

may not be regular. Exact typing may be done in the form of backward type inference, where the output type is given and a largest input type is inferred [20]. It would be interesting to use the backward approach to type our accumulators without the approximation introduced for “cons” for instance.

For what concerns XQuery and XPath, several complementary works are of interest. First, the work of Genevès *et al.* which encodes XPath and XQuery in the μ -calculus ([14,15] where zippers to manage XPath reverse axes were first introduced) supports our claim. Adding path expressions at the level of *types* is not more expensive: subtyping (or equivalently satisfiability of particular formulæ of the μ -calculus which are equivalent to regular tree languages) remains EXPTIME, even with upward paths (or in our case, zipper types). In contrast, typing path expressions and more generally XQuery programs is still a challenging topic. While the W3C’s formal semantics of XQuery [24] gives a polynomial time typechecking algorithm for XQuery (in the absence of nested “let” or “for” constructs), it remains far too imprecise (in particular, reverse axes are left untyped). Recently, Genevès *et al.* [13] also studied a problem of typing reverse axes by using regular expressions of μ -calculus formulæ as types, which they call focused-tree types. Since, as our zipped types, focused-tree types can describe both the type of the current node and its context, their type system also gives a precise type for reverse axis expressions. However, while focused-tree types are more concise than zipper types, it is difficult to type construction of a new XML document, and thus their type system requires an explicit type annotation for each XML element. Furthermore, their type system does not feature arrow types. That said, it will be quite interesting to combine their approach with ours.

We are currently implementing axis patterns and XPath expressions on top of the CDuce compiler. Future work includes extensions to other XQuery constructs as well as XMLSchema, the addition of aggregate functions by associating accumulators to specific operators, the inclusion of navigational expressions in types so as to exploit the full expressivity of our zipped types (*e.g.*, to type functions that work on the ancestors of their arguments), and the application of the polymorphic type system of [5,6] to both XQuery and navigational CDuce so that for instance the function *pretty* defined in the introduction can be given the following, far more precise intersection of two arrow types:

```
(<a class="style1" href= $\beta$  ..> $\gamma$  -> <a href= $\beta$ >[<b> $\gamma$ ])
& ( $\alpha$ \<a class="style1" href=_ ..>_ ->  $\alpha$ \<a class="style1" href=_ ..>_)
```

This type (where α , β , and γ denote universally quantified type variables) precisely describes, by the arrow type on the first line, the transformation of the sought links, and states, by the arrow on the second line, that in all the other cases (*i.e.*, for every type α different from the sought link) it returns the same type as the input. This must be compared with the corresponding type in Figure 1, where the types of the attribute href, of the content of the a element, and above all of any other value not matched by the first branch are not preserved.

Acknowledgments. We want to thank the reviewers of ESOP who gave detailed suggestions to improve our presentation. This work was partially supported by the ANR TYPEX project n. ANR-11-BS02-007.

References

1. N. Alon, T. Milo, F. Neven, D. Suciu, and V. Vianu. XML with data values: typechecking revisited. In *PODS*, pages 138–149. ACM, 2001.
2. R. M. Amadio and L. Cardelli. Subtyping recursive types. *ACM Trans. Program. Lang. Syst.*, 15(4):575–631, 1993.
3. M. Benedikt and H. Vu. Higher-order functions and structured datatypes. In *WebDB*, pages 43–48, 2012.
4. V. Benzaken, G. Castagna, and A. Frisch. CDuce: An XML-centric general-purpose language. In *ICFP*, pages 51–63, 2003.
5. G. Castagna, K. Nguyễn, Z. Xu, and P. Abate. Polymorphic functions with set-theoretic types. Part 2: Local type inference and type reconstruction. In *POPL*, pages 289–302, 2015.
6. G. Castagna, K. Nguyễn, Z. Xu, H. Im, S. Lenglet, and L. Padovani. Polymorphic functions with set-theoretic types. Part 1: Syntax, semantics, and evaluation. In *POPL*, pages 5–17, 2014.
7. G. Castagna and Z. Xu. Set-theoretic foundation of parametric polymorphism and subtyping. In *ICFP*, pages 94–106, 2011.
8. J. Engelfriet and H. Vogler. Macro tree transducers. *J. Comput. Syst. Sci.*, 31(1):71–146, 1985.
9. A. Frisch. *Théorie, conception et réalisation d’un langage adapté à XML*. PhD thesis, Université Paris 7 Denis Diderot, 2004.
10. A. Frisch, G. Castagna, and V. Benzaken. Semantic subtyping: Dealing set-theoretically with function, union, intersection, and negation types. *J. ACM*, 55(4):1–64, 2008.
11. V. Gapeyev, F. Garillot, and B. C. Pierce. Statically typed document transformation: An Xtatic experience. In *PLAN-X*, 2006.
12. V. Gapeyev and B. C. Pierce. Paths into patterns. Technical Report MS-CIS-04-25, University of Pennsylvania, Oct. 2004.
13. P. Genevès, N. Gesbert, and N. Layaïda. Xquery and static typing: Tackling the problem of backward axes. Available at <http://hal.inria.fr/hal-00872426>, July 2014.
14. P. Genevès and N. Layaïda. Eliminating dead-code from XQuery programs. In *ICSE*, 2010.
15. P. Genevès, N. Layaïda, and A. Schmitt. Efficient static analysis of XML paths and types. In *PLDI*, 2007.
16. H. Hosoya and B. C. Pierce. Regular expression pattern matching for XML. *J. Funct. Program.*, 13(6):961–1004, 2003.
17. H. Hosoya and B. C. Pierce. XDuce: A statically typed XML processing language. *ACM Trans. Internet Technol.*, 3(2):117–148, 2003.
18. G. Huet. The Zipper. *J. Funct. Program.*, 7(5):549–554, 1997.
19. R. Lämmel. Scrap your boilerplate with XPath-like combinators. In *POPL*, 2007.
20. S. Maneth, A. Berlea, T. Perst, and H. Seidl. XML type checking with macro tree transducers. In *PODS*, 2005.
21. W3C: XPath 1.0. <http://www.w3.org/TR/xpath>, 1999.
22. W3C: XPath 2.0. <http://www.w3.org/TR/xpath20>, 2010.
23. W3C: XML Query. <http://www.w3.org/TR/xquery>, 2010.
24. XQuery 1.0 and XPath 2.0 Formal Semantics (Second Edition). <http://www.w3.org/TR/xquery-semantics/>, 2010.
25. W3C: XQuery 3.0. <http://www.w3.org/TR/xquery-3.0>, 2014.
26. W3C: XML Schema. <http://www.w3.org/XML/Schema>, 2009.

Appendix

A Subtyping

In order to formally extend the subtyping relation defined in [10] to the types of Section 2.2, it suffices to modify DEFINITION 4.3 of [10] as done by Definition 15 (the reader can refer to [10] for the complete definitions of the notations used there). However, before extending the subtyping relation, we need some auxiliary definitions to give an interpretation of zipper types that is compatible with their infinite nature.

Let \mathcal{T}^Z denote the set of zipper types. First of all, notice that the contractiveness condition on zipper types implies that the binary relation $\triangleright \subseteq \mathcal{T}^Z \times \mathcal{T}^Z$ defined by $\tau_1 \vee \tau_2 \triangleright \tau_i$ and $\neg\tau \triangleright \tau$ is Noetherian (that is, strongly normalizing). This gives an induction principle on \mathcal{T}^Z that we use below.

Let D be a set and $\llbracket _ \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(D)$ be a set-theoretic interpretation (as defined in Definition 4.1 in [10], where \mathcal{T} denotes the set of all types and $\mathcal{P}(D)$ is the powerset of D). Let \mathcal{Z} denote $\{\text{L}d \mid d \in D\} \cup \{\text{R}d \mid d \in D\}$. We use \mathcal{Z}^* to denote the free monoid on \mathcal{Z} .

We define a binary predicate $(s \in_{\llbracket _ \rrbracket} \tau)$ that is parametric in the set-theoretic interpretation $\llbracket _ \rrbracket$ where $s \in \mathcal{Z}^*$ and $\tau \in \mathcal{T}^Z$. The truth value of $(s \in_{\llbracket _ \rrbracket} \tau)$ is defined by induction on the pair (s, τ) ordered lexicographically, using the inductive structure for elements of \mathcal{Z}^* (these are *finite* sequences s of decorated elements of D) and the induction principle we mentioned above for zipper types. Here is the definition:

$$\begin{aligned}
s \in_{\llbracket _ \rrbracket} \top &= \text{true} \\
\epsilon \in_{\llbracket _ \rrbracket} \bullet &= \text{true} \\
\text{R}d \cdot s \in_{\llbracket _ \rrbracket} \text{R}(u)_{\tau} \cdot \tau &= (s \in_{\llbracket _ \rrbracket} \tau) \text{ and } (d \in \llbracket u \rrbracket) \\
\text{L}d \cdot s \in_{\llbracket _ \rrbracket} \text{L}(u)_{\tau} \cdot \tau &= (s \in_{\llbracket _ \rrbracket} \tau) \text{ and } (d \in \llbracket u \rrbracket) \\
s \in_{\llbracket _ \rrbracket} \tau_1 \vee \tau_2 &= (s \in_{\llbracket _ \rrbracket} \tau_1) \text{ or } (s \in_{\llbracket _ \rrbracket} \tau_2) \\
s \in_{\llbracket _ \rrbracket} \neg\tau &= \text{not}(s \in_{\llbracket _ \rrbracket} \tau) \\
s \in_{\llbracket _ \rrbracket} \tau &= \text{false} \quad \text{otherwise}
\end{aligned}$$

This predicate is then used to define the following interpretation of zipper types:

Definition 14 (Zipper type interpretation). *Let D be a set, $\llbracket _ \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(D)$ be a set-theoretic interpretation, and \mathcal{Z} denote $\{\text{L}d \mid d \in D\} \cup \{\text{R}d \mid d \in D\}$. The interpretation $\llbracket \tau \rrbracket$ of a zipper type τ with respect to $\llbracket _ \rrbracket$ is defined as:*

$$\llbracket \tau \rrbracket = \{s \in \mathcal{Z}^* \mid s \in_{\llbracket _ \rrbracket} \tau\}$$

Subtyping for zipper types is defined as $\tau <: \tau' \stackrel{\text{def}}{=} \llbracket \tau \rrbracket \subseteq \llbracket \tau' \rrbracket$.

Finally, the interpretation of zipper types is used to extend DEFINITION 4.3 of [10] to our new types, as follows.

Definition 15 (Extensional interpretation). *Let $\llbracket _ \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(D)$ be a set-theoretic interpretation in some set D . Let $\mathcal{Z} \stackrel{\text{def}}{=} \{\text{L}d \mid d \in D\} \cup \{\text{R}d \mid d \in D\}$ and \mathcal{Z}^* denote the free monoid on \mathcal{Z} .*

We define the associated extensional interpretation as the unique set-theoretic interpretation

$$\mathbb{E}(_) : \mathcal{T} \rightarrow \mathcal{P}(\mathbb{E}D)$$

(where $\mathbb{E}D = \mathcal{C} + D^2 + \mathcal{P}(D \times D_\Omega) + (\mathcal{P}(D) \times \mathcal{Z}^*)$) such that:

$$\begin{aligned} \mathbb{E}(b) &= \mathbb{B}[\![b]\!] && \subseteq \mathcal{C} \\ \mathbb{E}(t_1 \times t_2) &= \mathbb{[[}t_1\]] \times \mathbb{[[}t_2\]] && \subseteq D^2 \\ \mathbb{E}(t_1 \rightarrow t_2) &= \mathbb{[[}t_1\]] \rightarrow \mathbb{[[}t_2\]] && \subseteq \mathcal{P}(D \times D_\Omega) \\ \mathbb{E}((u)_\tau) &= \mathbb{[[}u\]] \times \mathbb{[[}\tau\]] && \subseteq \mathcal{P}(D) \times \mathcal{Z}^* \end{aligned}$$

All the other definitions of [10] remain unchanged, in particular, those of a well-founded model and of its induced subtyping relation.

In order to decide the subtyping relation induced by a model, a possibility is to extend the definitions of Section 6 in [10] to account for the new zipper type constructor. A simpler way is to use a well-founded model, encode both zipper types and zipper values in the types and values of [10] (our pre-type and pre-values) via the encoding function $\text{Enc}(_)$ below, and prove that $\mathbb{[[}\tau_1\]] \subseteq \mathbb{[[}\tau_2\]]$ if and only if $\text{Enc}(\tau_1) \leq \text{Enc}(\tau_2)$.

Definition 16 (Encoding of zippers). *Zippers and zipper types are encoded (inductively) into pairs and (coinductively) into product types, respectively, as follows:*

$$\begin{aligned} \text{Enc}(\mathbb{L}(w)_\delta \cdot \delta) &\equiv ((\mathbb{L}, w), \text{Enc}(\delta)) \\ \text{Enc}(\mathbb{R}(w)_\delta \cdot \delta) &\equiv ((\mathbb{R}, w), \text{Enc}(\delta)) \\ \text{Enc}(\bullet) &\equiv \text{'nil} \\ \\ \text{Enc}(\mathbb{L}(u)_\tau \cdot \tau) &\equiv (\mathbb{L} \times u) \times \text{Enc}(\tau) \\ \text{Enc}(\mathbb{R}(u)_\tau \cdot \tau) &\equiv (\mathbb{R} \times u) \times \text{Enc}(\tau) \\ \text{Enc}(\top) &\equiv \mu X.((\mathbb{L} \vee \mathbb{R}) \times \mathbb{1}) \times X \vee \text{'nil} \\ \text{Enc}(\bullet) &\equiv \text{'nil} \\ \text{Enc}(\neg\tau) &\equiv \text{Enc}(\top) \setminus \text{Enc}(\tau) \\ \text{Enc}(\tau_1 \vee \tau_2) &\equiv \text{Enc}(\tau_1) \vee \text{Enc}(\tau_2) \end{aligned}$$

The termination of the subtyping algorithm in [10] implies the termination of the subtyping algorithm on (the encoding of) our extended type algebra. Since the encoding is linear on the size of terms, both algorithms have the same complexity.

As an aside, note that the whole calculus presented in this paper can be faithfully encoded in CDuce without affecting the complexity of the algorithms: it is straightforward to extend the encodings of Definition 16 to \mathcal{E} and \mathcal{T} .

All it remains to prove is the soundness and completeness of the encoding, namely:

Theorem 17. *Let \leq be a subtyping relation for the calculus in [10] induced by a well-founded model. Then, there exist a set D and a set-theoretic interpretation $\mathbb{[[}_]\!] : \mathcal{T} \rightarrow \mathcal{P}(D)$ such that for all $\tau_1, \tau_2 \in \mathcal{T}^Z$, the following holds:*

$$\mathbb{[[}\tau_1\]] \subseteq \mathbb{[[}\tau_2\]] \iff \text{Enc}(\tau_1) \leq \text{Enc}(\tau_2)$$

Proof. Let us use λ_{FCB} to denote the λ -calculus defined in [10]. Let \leq be any subtyping relation for λ_{FCB} induced by a well-founded model. Theorem 5.5 in [10] states that $\text{Enc}(\tau_1) \leq \text{Enc}(\tau_2)$ if and only if $\llbracket \text{Enc}(\tau_1) \rrbracket_{\mathcal{V}} \subseteq \llbracket \text{Enc}(\tau_2) \rrbracket_{\mathcal{V}}$, where $\llbracket _ \rrbracket_{\mathcal{V}}$ is the *value interpretation* for the types of λ_{FCB} defined as $\llbracket t \rrbracket_{\mathcal{V}} \stackrel{\text{def}}{=} \{v \mid \vdash v : t\}$ (where v and t respectively range over the values and types of λ_{FCB}).

The simplest way to prove this theorem, then, is to produce a set D and interpretation $\llbracket _ \rrbracket : \mathcal{T} \rightarrow \mathcal{P}(D)$ such that there is a one-to-one correspondence between $\llbracket \tau \rrbracket$ and $\llbracket \text{Enc}(\tau) \rrbracket_{\mathcal{V}}$.

To that end, take as D the set of all values of λ_{FCB} , that is, $\llbracket \mathbb{1} \rrbracket_{\mathcal{V}}$, and as interpretation any interpretation that on the pre-values of our calculus (which are the values of λ_{FCB}) behaves as $\llbracket _ \rrbracket_{\mathcal{V}}$, that is, $\llbracket w \rrbracket = \llbracket w \rrbracket_{\mathcal{V}}$ for every pre-value w . Next, we define an embedding function $f : \mathcal{Z}^* \hookrightarrow \llbracket \text{Enc}(\top) \rrbracket_{\mathcal{V}}$ from the resulting \mathcal{Z}^* to set of values of type $\text{Enc}(\top)$, by induction on the length of the elements of \mathcal{Z}^* as follows:

$$\begin{aligned} f(\varepsilon) &= \text{'nil} \\ f(\text{L } w \cdot s) &= ((\text{'L}, w), f(s)) \\ f(\text{R } w \cdot s) &= ((\text{'R}, w), f(s)) \end{aligned}$$

We can prove that f is injective by induction on \mathcal{Z}^* and surjective by induction on $\llbracket \text{Enc}(\top) \rrbracket_{\mathcal{V}}$ (recall that, contrary to types, values are inductively defined).

Since $s \in \llbracket \tau \rrbracket \iff s \in_{\llbracket \llbracket _ \rrbracket \rrbracket} \tau$, then to prove that f is a one-to-one mapping from $\llbracket \tau \rrbracket$ to $\llbracket \text{Enc}(\tau) \rrbracket_{\mathcal{V}}$, it suffices to prove that for all $s \in \mathcal{Z}^*$ and $\tau \in \mathcal{T}^Z$, $s \in_{\llbracket \llbracket _ \rrbracket \rrbracket} \tau \iff f(s) \in \llbracket \text{Enc}(\tau) \rrbracket_{\mathcal{V}}$. This can be easily proved by induction on the pair (s, τ) ordered lexicographically, by performing a case analysis on the definition of $\in_{\llbracket \llbracket _ \rrbracket \rrbracket}$.

Zipper typing rules

 $\boxed{\vdash \delta : \tau}$

$$\begin{array}{c} \text{[ZT-ROOT]} \\ \hline \vdash \bullet : \bullet \end{array} \quad \begin{array}{c} \text{[ZT-SUB]} \\ \frac{\vdash \delta : \tau \quad \tau <: \tau'}{\vdash \delta : \tau'} \end{array} \quad \begin{array}{c} \text{[ZT-LEFT]} \\ \frac{\vdash w : t \quad \vdash \delta : \tau}{\vdash \mathbf{L}(w)_\delta \cdot \delta : \mathbf{L}(t)_\tau \cdot \tau} \end{array} \quad \begin{array}{c} \text{[ZT-RIGHT]} \\ \frac{\vdash w : t \quad \vdash \delta : \tau}{\vdash \mathbf{R}(w)_\delta \cdot \delta : \mathbf{R}(t)_\tau \cdot \tau} \end{array}$$

Typing rules

 $\boxed{\Gamma \vdash e : t}$

$$\begin{array}{c} \text{[T-CST]} \quad \text{[T-VAR]} \quad \text{[T-ACC]} \quad \text{[T-ZIP-VAL]} \\ \hline \Gamma \vdash c : b_c \quad \Gamma \vdash x : \Gamma(x) \quad \Gamma \vdash \hat{x} : \Gamma(\hat{x}) \quad \frac{\vdash w : t \quad \vdash \delta : \tau \quad t \leq \mathbb{1}_{\text{NZ}}}{\Gamma \vdash (w)_\delta : (t)_\tau} \end{array}$$

$$\begin{array}{c} \text{[T-ZIP-EXPR]} \quad \text{[T-PAIR]} \quad \text{[T-SUB]} \\ \frac{\vdash e : t \quad t \leq \mathbb{1}_{\text{NZ}}}{\Gamma \vdash (e)_\bullet : (t)_\bullet} \quad \frac{\Gamma \vdash e_1 : t_1 \quad \Gamma \vdash e_2 : t_2}{\Gamma \vdash (e_1, e_2) : t_1 \times t_2} \quad \frac{\Gamma \vdash e : s \quad s \leq t}{\Gamma \vdash e : t} \end{array}$$

$$\begin{array}{c} \text{[T-OP]} \\ \frac{\forall i = 1..n_o, \Gamma \vdash e_i : t_i \quad t_1, \dots, t_{n_o} \xrightarrow{o} t}{\Gamma \vdash o(e_1, \dots, e_{n_o}) : t} \text{ for } o \in \mathcal{O} \end{array}$$

$$\begin{array}{c} \text{[T-FUN]} \\ \frac{t = \bigwedge_{i=1..n} (t_i \rightarrow s_i) \wedge \bigwedge_j \neg(t'_j \rightarrow s'_j) \quad t \not\equiv 0 \quad \forall i = 1..n, \Gamma \cup \{f \mapsto t, x \mapsto t_i\} \vdash e : s_i}{\Gamma \vdash \mu f^{(t_1 \rightarrow s_1; \dots; t_n \rightarrow s_n)}(x).e : t} \end{array}$$

$$\begin{array}{c} \text{[T-MATCH]} \\ \frac{\begin{array}{l} t \leq \wr p_1 \vee \wr p_2 \\ t_1 \equiv t \wedge \wr p_1 \quad t_2 \equiv t \wedge \neg \wr p_1 \\ \Sigma_i \equiv \{\hat{x} \mapsto \text{Init}(\hat{x}) \mid \hat{x} \in \text{Acc}(p_i)\} \end{array} \quad \begin{array}{l} \Gamma \vdash e : t \\ \Gamma_i \equiv \square \vdash t_i / p_i \quad \Gamma'_i \equiv \Sigma_i; \square \vdash t_i // p_i \\ \Gamma \cup \Gamma_i \cup \Gamma'_i \vdash e_i : t'_i \end{array}}{\Gamma \vdash \text{match } e \text{ with } p_1 \rightarrow e_1 \mid p_2 \rightarrow e_2 : \bigvee_{\{i \mid t_i \neq 0\}} t'_i} \quad (i = 1, 2) \end{array}$$

Fig. 10: Typing rules

B Type System

B.1 Typing Rules

Figure 10 shows typing rules for our calculus in the form of judgments $\vdash \delta : \tau$ and $\Gamma \vdash e : t$ where the former is for typing zippers and the latter for typing expressions.

Theorem 18 (Accepted types). *Let p be a pattern such that for every \hat{x} in $\text{Acc}(p)$, $\text{Op}(\hat{x})$ has an exact input. Then, the set of all values v such that*

$$\{\hat{x} \mapsto \text{Init}(\hat{x}) \mid \hat{x} \in \text{Acc}(p)\}; \square \vdash v/p \not\prec \Omega$$

is a type. We call this set the accepted type of p and denote it by $\wr p$. It can be computed by solving the following guarded system of equations (the variables are

respectively the $\wr p'$ and $\wr \varphi$ for the subterms p' and φ of p).

$$\begin{array}{ll}
\wr t \wr & = t & \wr \tau \wr & = \tau \\
\wr x \wr & = \mathbb{1} & \wr \mathbf{L} p \cdot \varphi \wr & = \mathbf{L} \wr p \wr \cdot \wr \varphi \wr \\
\wr \dot{x} \wr & = \mathbb{I}(\mathbf{Op}(\dot{x})) & \wr \mathbf{R} p \cdot \varphi \wr & = \mathbf{R} \wr p \wr \cdot \wr \varphi \wr \\
\wr (p_1, p_2) \wr & = \wr p_1 \wr \times \wr p_2 \wr & \wr \varphi_1 \mid \varphi_2 \wr & = \wr \varphi_1 \wr \vee \wr \varphi_2 \wr \\
\wr p_1 \mid p_2 \wr & = \wr p_1 \wr \vee \wr p_2 \wr \\
\wr p_1 \ \& p_2 \wr & = \wr p_1 \wr \wedge \wr p_2 \wr \\
\wr (x := c) \wr & = \mathbb{1} \\
\wr (q)_\varphi \wr & = (\wr q \wr)_{\wr \varphi \wr}
\end{array}$$

Below, we introduce a technical notation for the matching of product types. Indeed, the most general type for any pair is a *finite* union of products (or, said differently, while it is possible to push intersections below product constructors, it is not possible to do it for unions without introducing an approximation since in general: $(a \times b) \vee (c \times d) \preceq (a \vee c) \times (b \vee d)$).

Lemma 19 (Product decomposition). *Let t be a type such that $t \leq \mathbb{1} \times \mathbb{1}$. Then there exists a finite set of pairs of types*

$$\Pi(t) = \bigcup_{i \leq n} \{(t_1^i, t_2^i)\}$$

such that

$$t \simeq \bigvee_{(t_1, t_2) \in \Pi(t)} t_1 \times t_2$$

Furthermore, given a decomposition $\Pi(t)$, we define the first and second type projection as:

$$\Pi_i(t) = \bigcup_{(t_1, t_2) \in \Pi(t)} \{t_i\}$$

There exist many such product decompositions. For instance, one decomposition is obtained by taking the syntactic expression given for a product type t and pushing intersections below products until only unions remain at top-level. More complex decompositions (which yield more precise typing or more efficient pattern matching) are described in [9]. As zipped types $(u)_\tau$ and zipper types of the form $\mathbf{L} t \cdot \tau$ and $\mathbf{R} t \cdot \tau$ are essentially product types, we also introduce similar decompositions for such types.

Lemma 20 (Zipped type decomposition). *Let t be a type such that $t \leq (\mathbb{1})_\top$. Then there exists a finite set of pairs of types*

$$\mathcal{Z}(t) = \bigcup_{i \leq n} \{(u_i, \tau_i)\}$$

such that

$$t \simeq \bigvee_{(u, \tau) \in \mathcal{Z}(t)} (u)_\tau$$

Furthermore, given a decomposition $\mathcal{Z}(t)$, we define the first and second type projection as

$$\mathcal{Z}_1(t) = \bigcup_{(u,\tau) \in \mathcal{Z}(t)} \{u\} \quad \mathcal{Z}_2(t) = \bigcup_{(u,\tau) \in \mathcal{Z}(t)} \{\tau\}$$

Note that zipped types are just a special form of product types and therefore we can use the same product decomposition method for zipped typed decomposition.

We also need another product decomposition called *zipper type decomposition* for zipper types $\bigvee_{i \leq n} \mathsf{L}(u_i)_{\tau_i} \cdot \tau_i$ and $\bigvee_{j \leq m} \mathsf{R}(u'_j)_{\tau'_j} \cdot \tau'_j$. Note that they are also a special form of product types.

Lemma 21 (Zipper type decomposition). *Let τ be a zipper type such that $\tau \leq \mathsf{L}\mathbb{1} \cdot \top$. Then there exists a finite set of pairs of types and zipper types*

$$\mathcal{Z}^{\mathsf{L}}(\tau) = \bigcup_{i \leq n} \{(t_i, \tau_i)\}$$

such that

$$\tau \simeq \bigvee_{(t_i, \tau_i) \in \mathcal{Z}^{\mathsf{L}}(\tau)} \mathsf{L} t_i \cdot \tau_i$$

As before, we define the first and second type projection of $\mathcal{Z}^{\mathsf{L}}(\tau)$ as

$$\mathcal{Z}_1^{\mathsf{L}}(\tau) = \bigcup_{(t, \tau') \in \mathcal{Z}^{\mathsf{L}}(\tau)} \{t\} \quad \mathcal{Z}_2^{\mathsf{L}}(\tau) = \bigcup_{(t, \tau') \in \mathcal{Z}^{\mathsf{L}}(\tau)} \{\tau'\}$$

We similarly define decomposition operations $\mathcal{Z}^{\mathsf{R}}(\tau)$ and $\mathcal{Z}_i^{\mathsf{R}}(\tau)$ for a zipper type τ such that $\tau \leq \mathsf{R}\mathbb{1} \cdot \top$.

Lemma 22 (Optional zipper types). *We write $\tau^?$ to denote an optional zipper type, which is either τ or a none zipper type \square .*

$$\tau^? ::= \tau \mid \square$$

For notational simplicity, we consider $(t)_{\square}$ to be equivalent to t .

B.2 Type Environment for Capture Variables

Figure 11 defines a set of recursive equations of the forms $(\tau^? \vdash t/p)(x) = t'$ and $(\tau/\varphi)(y) = t''$ with the following implicit assumptions: $t \leq \llbracket p \rrbracket$, $x \in \text{Var}(p)$, $\tau \leq \llbracket \varphi \rrbracket$, and $y \in \text{Var}(\varphi)$. Basically, given a type t , a pattern p , and an optional zipper type $\tau^?$, $\tau^? \vdash t/p$ computes a type environment for capture variables used in the pattern p and similarly for τ/φ . In the figure, in order to properly update zipper types for navigational patterns (*i.e.*, pair patterns), we use the following additional projection operations.

For types:

$$\begin{aligned}
(\tau^? \vdash t/x)(x) &= (t)_{\tau^?} \\
(\tau^? \vdash t/(p_1, p_2))(x) &= \bigvee_{(t_1, t_2) \in \Pi(t)} (\Pi_1(t_1 \times t_2, \tau^?) \vdash t_1/p_1)(x) && \text{if } x \in \text{Var}(p_1) \setminus \text{Var}(p_2) \\
(\tau^? \vdash t/(p_1, p_2))(x) &= \bigvee_{(t_1, t_2) \in \Pi(t)} (\Pi_2(t_1 \times t_2, \tau^?) \vdash t_2/p_2)(x) && \text{if } x \in \text{Var}(p_2) \setminus \text{Var}(p_1) \\
(\tau^? \vdash t/(p_1, p_2))(x) &= \bigvee_{(t_1, t_2) \in \Pi(t)} (\Pi_1(t_1 \times t_2, \tau^?) \vdash t_1/p_1)(x) \times (\Pi_2(t_1 \times t_2, \tau^?) \vdash t_2/p_2)(x) && \text{if } x \in \text{Var}(p_1) \cap \text{Var}(p_2) \\
(\tau^? \vdash t/p_1 \mid p_2)(x) &= (\tau^? \vdash (t \wedge \lambda p_1 \mathbin{\text{f}})/p_1)(x) \vee (\tau^? \vdash (t \wedge \neg \lambda p_1 \mathbin{\text{f}})/p_2)(x) \\
(\tau^? \vdash t/p_1 \& p_2)(x) &= (\tau^? \vdash t/p_i)(x) && \text{if } x \in \text{Var}(p_i) \\
(\tau^? \vdash t/(x := c))(x) &= t_c && \text{if } t \not\approx \emptyset \\
(\tau^? \vdash t/(x := c))(x) &= \emptyset && \text{if } t \simeq \emptyset \\
(\square \vdash t/(q)_\varphi)(x) &= \bigvee_{(u, \tau) \in \mathcal{Z}(t)} (\tau \vdash u/q)(x) && \text{if } x \in \text{Var}(p) \\
(\square \vdash t/(q)_\varphi)(x) &= (\mathcal{Z}_2(t)/\varphi)(x) && \text{if } x \in \text{Var}(\varphi)
\end{aligned}$$

For zipper types:

$$\begin{aligned}
(\tau/Lp \cdot \varphi)(x) &= (\square \vdash \mathcal{Z}_1^L(\tau)/p)(x) && \text{if } x \in \text{Var}(p) \setminus \text{Var}(\varphi) \\
(\tau/Lp \cdot \varphi)(x) &= (\mathcal{Z}_2^L(\tau)/\varphi)(x) && \text{if } x \in \text{Var}(\varphi) \setminus \text{Var}(p) \\
(\tau/Lp \cdot \varphi)(x) &= \bigvee_{(t, \tau') \in \mathcal{Z}^L(\tau)} ((\square \vdash t/p)(x)) \times ((\tau'/\varphi)(x)) && \text{if } x \in \text{Var}(p) \cap \text{Var}(\varphi) \\
(\tau/Rp \cdot \varphi)(x) &= (\square \vdash \mathcal{Z}_1^R(\tau)/p)(x) && \text{if } x \in \text{Var}(p) \setminus \text{Var}(\varphi) \\
(\tau/Rp \cdot \varphi)(x) &= (\mathcal{Z}_2^R(\tau)/\varphi)(x) && \text{if } x \in \text{Var}(\varphi) \setminus \text{Var}(p) \\
(\tau/Rp \cdot \varphi)(x) &= \bigvee_{(t, \tau') \in \mathcal{Z}^R(\tau)} ((\square \vdash t/p)(x)) \times ((\tau'/\varphi)(x)) && \text{if } x \in \text{Var}(p) \cap \text{Var}(\varphi) \\
(\tau/\varphi_1 \mid \varphi_2)(x) &= ((\tau \wedge \lambda \varphi_1 \mathbin{\text{f}})/\varphi_1)(x) \vee ((\tau \wedge \neg \lambda \varphi_1 \mathbin{\text{f}})/\varphi_2)(x)
\end{aligned}$$

Fig. 11: Computing the type environment for capture variables

Definition 23. Let t be a type and $\tau^?$ an optional zipper type. Then, $\Pi_i(t, \tau^?)$ computes the left projection of the zipped type $(t)_{\tau^?}$ when $i = 1$ and the right projection when $i = 2$ as follows:

$$\begin{aligned} \Pi_1(u, \tau) &\stackrel{\text{def}}{=} \mathbf{L}(u)_{\tau} \cdot \tau \\ \Pi_2(u, \tau) &\stackrel{\text{def}}{=} \mathbf{R}(u)_{\tau} \cdot \tau \\ \Pi_i(t, \square) &\stackrel{\text{def}}{=} \square \end{aligned}$$

Although the equations in Figure 11 seem to be much more complicated than those in [9], if we ignore the zipper-related part, they are exactly the same. For example, when $\tau^? = \square$ and $x \in \text{Var}(p_1) \setminus \text{Var}(p_2)$, the equation $(\tau^? \vdash t/(p_1, p_2))(x)$ amounts to $\bigvee_{(t_1, t_2) \in \Pi(t)} (t_1/p_1)(x)$, which is equal to $(\Pi_1(t)/p_1)(x)$ as in [9]. Of course, even in the presence of zipper types, we can have a simpler equation such as $(\tau^? \vdash t/(p_1, p_2))(x) = (\Pi_1(t, \tau^?) \vdash \Pi_1(t)/p_1)(x)$ at the expense of more precise typing.

Unfortunately, the set of equations in Figure 11 is not well-founded. The main source of the problem is the equations for pair patterns and the use of the projection function $\Pi_i(_, _)$ defined in Definition 23. More precisely, given a zipped recursive type and a zipped recursive pattern, those equations may in general generate an infinite number of new equations. For instance, consider an integer list type, defined as a recursive equation $X = (\text{int} \times X) \vee \text{nil}$, a recursive pattern $Y = (y, Y) \mid (y := \text{nil})$, and a matching of a zipped type $(X)_{\top}$ against a zipped pattern $(Y)_{\top}$, that is, $\square \vdash (X)_{\top}/(Y)_{\top}$. Then, the matching generates the following infinite set of equations:

$$\begin{aligned} (\square \vdash (X)_{\top}/(Y)_{\top})(y) &= (\top \vdash X/Y)(y) \\ &= (\top \vdash (\text{int} \times X)/(y, Y))(y) \vee \text{nil} \\ &= ((\Pi_1(\text{int} \times X, \top) \vdash \text{int}/y)(y) \times (\Pi_2(\text{int} \times X, \top) \vdash X/Y)(y)) \vee \text{nil} \\ &= ((\text{int})_{\mathbf{L}(\text{int} \times X)_{\top} \cdot \top} \times (\mathbf{R}(\text{int} \times X)_{\top} \cdot \top \vdash X/Y)(y)) \vee \text{nil} \\ &= (\dots \times ((\mathbf{R}(\text{int} \times X)_{\top} \cdot \top \vdash \text{int} \times X/(y, Y))(y)) \vee \text{nil}) \vee \text{nil} \\ &= \dots \end{aligned}$$

The intuition behind this infinite sequence of equations is as follows. Given a sequence of integers, zipper annotations specify for each integer its order in the sequence. Furthermore, the integer list type (*i.e.*, X in the above equations) represents infinitely many sequences of integers of an arbitrary length and for each sequence of a fixed length there is the unique most precise type with zipper annotations (one that the equations in Figure 11 try to compute). Therefore, the computation of $(\square \vdash (X)_{\top}/(Y)_{\top})(y)$ generates an infinite sequence of equations and does not terminate.

To remedy this problem, we simply introduce an approximation of $\Pi_i(_, _)$ in Definition 24, which guarantees the termination of the computation of a type environment. Note that in the above example equations, $\Pi_i(_, _)$ computes a new zipper type (by adding some prefix to the argument zipper type, that is, the second argument) which is used as an argument of $\Pi_i(_, _)$ in another equation, thus generating a new equation.

Definition 24. Given an optional zipper type $\tau^?$, we define an idempotent operation $\text{clos}(\tau^?)$ as follows:

$$\begin{aligned} \text{clos}(\Box) &\stackrel{\text{def}}{=} \Box \\ \text{clos}(\tau) &\stackrel{\text{def}}{=} \mu X.((L\mathbb{1} \cdot X) \vee (R\mathbb{1} \cdot X)) \vee \tau \quad (X \text{ fresh}) \end{aligned}$$

Intuitively $\text{clos}(\tau)$ adds a possibly infinite sequence of $L\mathbb{1}$ and $R\mathbb{1}$ to the given τ . Note that applying $\text{clos}(_)$ multiple times yields the same result as its initial application.

To ensure the termination of the computation of a type environment, we can use $\text{clos}(\tau^?)$ instead of $\Pi_i(t, \tau^?)$ in the equations for pair patterns in Figure 11 when both the input type and pattern are recursively defined (if either the input type or pattern is not recursive, the computation trivially terminates and the set of equations in Figure 11 computes a precise type for each capture variable). In practice, we may use $\Pi_i(_, _)$ up to a certain number of times and then from that point use $\text{clos}(_)$ to guarantee that the set of generated equations is finite.

B.3 Type Environment for Accumulators

Figure 12 defines a set of recursive equations of the forms $\Sigma; \tau^? \vdash t // p = \Sigma'$ and $\Sigma \vdash \tau // \varphi = \Sigma'$ with the following implicit assumptions: $t \leq \wr p$ and $\tau \leq \wr \varphi$. Basically, given a type t , a pattern p , an optional zipper type $\tau^?$, and an input type environment Σ for accumulators, $\Sigma; \tau^? \vdash t // p$ computes a type environment for accumulators by updating Σ (for example, see the third equation in the figure) and similarly for $\Sigma \vdash \tau // \varphi$. Due to the same reason as the equations in Figure 11, the set of equations in Figure 12 is not well-founded, and in order to ensure the termination of the computation of a type environment, we use the same technique as before: we use $\text{clos}(_)$ instead of $\Pi_i(_, _)$ in the equations.

B.4 Zipper Erasure

The top-level erasure simply removes the top-level zipper type annotation, while the deep erasure is typed by recursively removing the zipper annotations from the input type as follows:

$$\begin{aligned} (t)_\tau &\xrightarrow{\text{rm}} t && \text{if } t \wedge (\mathbb{1})_\top \simeq \mathbb{0} \\ t &\xrightarrow{\text{rm}} (t \wedge \neg(\mathbb{1})_\top) \vee s && \text{where } t \wedge (\mathbb{1})_\top \xrightarrow{\text{rm}} s \\ \\ t &\xrightarrow{\text{dr}} t && \text{if } t \leq \mathbb{1}_{\text{NZ}} \\ t &\xrightarrow{\text{dr}} t \wedge (\mathbb{1}_{\text{basic}} \vee \mathbb{1}_{\text{fun}}) && \\ &\vee \bigvee_{(t_1, t_2) \in \Pi(t \wedge \mathbb{1}_{\text{prod}})} t'_1 \times t'_2 && \text{where } t_i \xrightarrow{\text{dr}} t'_i \\ &\vee s && \text{where } t \wedge (\mathbb{1}_{\text{NZ}})_\top \xrightarrow{\text{rm}} s \end{aligned}$$

There are two cases for the deep erasure. If an input type does not contain any zipper type annotation (in-depth), it is left unchanged. Otherwise, the type is split into three parts. The first part consists of basic types and arrow types and is “copied” unchanged in the output type. The second part corresponds to the product type components of the union, and in this case the output type is the union of the products formed from the erasure of each component. The third part corresponds to zipped types in which the zipper type annotations are removed using the top-level erasure `rm`. We need to ensure that the `drm` function terminates, that is, given a type t , the number of t_i in the second part is finite and this property is shown in Section C).

For types:

$$\begin{aligned}
\Sigma; \tau^2 \vdash t // t' &= \Sigma \\
\Sigma; \tau^2 \vdash t // x &= \Sigma \\
\Sigma; \tau^2 \vdash t // \dot{x} &= \Sigma[s/\dot{x}] && \text{if } (t)_{\tau^?}, \Sigma(\dot{x}) \xrightarrow{\text{Op}(\dot{x})} s \\
\Sigma; \tau^2 \vdash t // (p_1, p_2) &= \bigsqcup_{(t_1, t_2) \in \Pi(t)} (\Sigma; \Pi_1(t_1 \times t_2, \tau^2) \vdash t_1 // p_1; \Pi_2(t_1 \times t_2, \tau^2) \vdash t_2 // p_2) \\
\Sigma; \tau^2 \vdash t // p_1 \mid p_2 &= \Sigma; \tau^2 \vdash t // p_1 && \text{if } t \leq \wr p_1 \wr \\
\Sigma; \tau^2 \vdash t // p_1 \mid p_2 &= \Sigma; \tau^2 \vdash t // p_2 && \text{if } t \leq \neg \wr p_1 \wr \\
\Sigma; \tau^2 \vdash t // p_1 \mid p_2 &= (\Sigma; \tau^2 \vdash (t \wedge \wr p_1 \wr) // p_1) \sqcup (\Sigma; \tau^2 \vdash (t \wedge \neg \wr p_1 \wr) // p_2) && \text{otherwise}
\end{aligned}$$

$$\Sigma; \tau^2 \vdash t // p_1 \& p_2 = (\Sigma; \tau^2 \vdash t // p_1); \tau^2 \vdash t // p_2$$

$$\Sigma; \tau^2 \vdash t // (x := c) = \Sigma$$

$$\Sigma; \square \vdash t // (q)_{\varphi} = \bigsqcup_{(u, \tau) \in \mathcal{Z}(t)} (\Sigma; \tau \vdash u // q) \vdash \tau // \varphi$$

For zipper types:

$$\begin{aligned}
\Sigma \vdash \tau // \tau' &= \Sigma \\
\Sigma \vdash \tau // \mathbf{L} p \cdot \varphi &= \bigsqcup_{(t, \tau') \in \mathcal{Z}^{\mathbf{L}}(\tau)} (\Sigma; \square \vdash t // p) \vdash \tau' // \varphi \\
\Sigma \vdash \tau // \mathbf{R} p \cdot \varphi &= \bigsqcup_{(t, \tau') \in \mathcal{Z}^{\mathbf{R}}(\tau)} (\Sigma; \square \vdash t // p) \vdash \tau' // \varphi \\
\Sigma \vdash \tau // \varphi_1 \mid \varphi_2 &= \Sigma \vdash \tau // \varphi_1 && \text{if } \tau \leq \wr \varphi_1 \wr \\
\Sigma \vdash \tau // \varphi_1 \mid \varphi_2 &= \Sigma \vdash \tau // \varphi_2 && \text{if } \tau \leq \neg \wr \varphi_1 \wr \\
\Sigma \vdash \tau // \varphi_1 \mid \varphi_2 &= (\Sigma \vdash (t \wedge \wr \varphi_1 \wr) // \varphi_1) \sqcup (\Sigma \vdash (t \wedge \neg \wr \varphi_1 \wr) // \varphi_2) && \text{otherwise}
\end{aligned}$$

$$(\Sigma_1 \sqcup \Sigma_2)(\dot{x}) = \begin{cases} \Sigma_1(\dot{x}) & \text{if } \dot{x} \in \text{dom}(\Sigma_1) \setminus \text{dom}(\Sigma_2) \\ \Sigma_2(\dot{x}) & \text{if } \dot{x} \in \text{dom}(\Sigma_2) \setminus \text{dom}(\Sigma_1) \\ \Sigma_1(\dot{x}) \vee \Sigma_2(\dot{x}) & \text{if } \dot{x} \in \text{dom}(\Sigma_1) \cap \text{dom}(\Sigma_2) \end{cases}$$

Fig. 12: Computing the type environment for accumulators

C Soundness Proofs

We first show that the algorithms we defined (such as the equations of Figure 12) or the typing of operator “ $\text{drm}(_)$ ” terminates, even for infinite (regular) input. To this end, we define the notion of plinth.

Definition 25 (Plinth). *Let \mathcal{O} be a set of operators. A plinth $\sqsupset_{\mathcal{O}} \subset \mathcal{T}$ over \mathcal{O} is a set of types with the following properties:*

Finiteness $\sqsupset_{\mathcal{O}}$ is finite;

Boolean closure $\sqsupset_{\mathcal{O}}$ contains $\mathbb{1}$ and $\mathbb{0}$ and is closed under Boolean connectives (\wedge, \vee, \neg) ;

Stability w.r.t. operators for all operators $o \in \mathcal{O}$ and types $t_1, \dots, t_{n_o} \in \sqsupset_{\mathcal{O}}$, if $t_1, \dots, t_{n_o} \xrightarrow{o} t$ then $t \in \sqsupset_{\mathcal{O}}$.

Intuitively, the plinth is the approximation of the set of types that can be found by saturating an initial set with the operators of \mathcal{O} . This means that any algorithm visiting types produced by the application of such operators will visit only a finite number of types and therefore terminate.

Lemma 26. *Let $\mathcal{O}_{\mathcal{Z}}$ be the set*

$$\mathcal{O}_{\mathcal{Z}} \equiv \{\mathcal{Z}_1, \mathcal{Z}_2, \mathcal{Z}_1^{\mathbb{L}}, \mathcal{Z}_2^{\mathbb{L}}, \mathcal{Z}_1^{\mathbb{R}}, \mathcal{Z}_2^{\mathbb{R}}\}$$

For all type t , there exists a plinth over $\mathcal{O}_{\mathcal{Z}}$ that contains t .

Proof. We assume here (see [9] for the proof) that given a finite set S of types, there exists a finite set $S' \supseteq S$ that is closed under \vee and \neg (and therefore also closed under \wedge which can be expressed in terms of union and negation).

We consider the initial set $S_0 = \{t, \mathbb{0}, \mathbb{1}, \bullet, \mathbb{L}\mathbb{1} \cdot \top, \mathbb{R}\mathbb{1} \cdot \mathbb{1}, \top\}$. Its saturation S'_0 with respect to \vee, \wedge, \neg is also finite. We now compute

$$S_1 = S'_0 \cup \bigcup_{t \in \{t \in S'_0 \mid t \leq (\mathbb{1})_{\top}\}} (\mathcal{Z}_1(t) \cup \mathcal{Z}_2(t))$$

By construction, S_1 is finite and so is its saturation S'_1 since zipper annotations cannot be nested. Then we compute

$$S_2 = S'_1 \cup \bigcup_{\tau \in \{\tau \in S'_1 \mid \tau \leq \mathbb{L}\mathbb{1} \cdot \top\}} (\mathcal{Z}_1^{\mathbb{L}}(\tau) \cup \mathcal{Z}_2^{\mathbb{L}}(\tau)) \cup \bigcup_{\tau \in \{\tau \in S'_1 \mid \tau \leq \mathbb{R}\mathbb{1} \cdot \top\}} (\mathcal{Z}_1^{\mathbb{R}}(\tau) \cup \mathcal{Z}_2^{\mathbb{R}}(\tau))$$

By construction, S_2 is finite and so is its saturation S'_2 due to the regularity condition on both types and zipper types. Finally, we compute S_3 by extending S'_2 using the \mathcal{Z}_1 operator, which is the plinth that we want. Note that \mathcal{Z}_1 is idempotent and that S'_2 is already saturated by \mathcal{Z}_2 because for any $t \in S'_2$, there are only two cases: if $t \in S'_0$ then $\mathcal{Z}_2(t) \subset S_1 \subset S'_2$; otherwise, $t \in S_2$ and it should be obtained from either $\mathbb{L}t \cdot \tau \in S'_1$ or $\mathbb{R}t \cdot \tau \in S'_1$, which implies the fact that t is of the form $(u)_{\tau}$ and τ is already included in S_2 and therefore S'_2 by either $\mathcal{Z}_2^{\mathbb{L}}(_)$ or $\mathcal{Z}_2^{\mathbb{R}}(_)$.

Lemma 27. *Let $\mathcal{O}_{\text{CDuce}}$ be the set*

$$\mathcal{O}_{\text{CDuce}} \equiv \mathcal{O}_{\mathcal{Z}} \cup \{\pi_1, \pi_2, \text{rm}, \text{cons}, \text{snoc}\}$$

For all type t , there exists a plinth over $\mathcal{O}_{\text{CDuce}}$ that contains t .

Proof. Let S_0 be a plinth over $\mathcal{O}_{\mathcal{Z}}$ that contains t , which can be computed by using Lemma 26. Then we compute

$$S_1 = S_0 \cup \bigcup_{t \in \{t \in S_0 \mid t \leq \mathbb{1}_{\text{prod}}\}} \pi_1(t)$$

By construction, S_1 is finite and so is its saturation S'_1 thanks to the regularity condition on types and zipper types. We similarly construct S'_2 (saturation of the closure by π_2) and S'_3 (saturation of the closure by rm). Next, consider the set $T \subseteq S'_3$ of types of the form $[(t_1 \vee \dots \vee t_n)*]$ that are valid for the second argument of cons and snoc . We compute

$$S_4 = S'_3 \cup \bigcup_{[(t_1 \vee \dots \vee t_n)*] \in T, U \subseteq S'_3} [((\bigvee_{s \in U} s) \vee t_1 \vee \dots \vee t_n)*]$$

This saturates S'_3 with respect to cons and snoc (which have the same typing function). Lastly, we compute the saturated set S'_4 , which is the plinth we seek. The crucial point here is that since $[(t_1 \vee \dots \vee t_n)*] \in S'_3$ which is in particular already saturated by π_1 and π_2 , $t_1 \vee \dots \vee t_n$ and 'nil' are already in S'_3 , and so is $\bigvee_{s \in U} s$ (since $U \subseteq S'_3$). Since S'_3 is saturated with respect to \vee , $(\bigvee_{s \in U} s) \vee t_1 \vee \dots \vee t_n \in S'_3$ and therefore S'_4 remains saturated by π_i ; and there is no need to iterate the process again.

Corollary 28. *The typing of operator drm terminates.*

Proof. Lemma 27 ensures that during the computation of drm , the operators π_1 , π_2 , and rm used in the definition of drm produce only a finite set of types. In other words, given a type t , drm needs to inspect only a finite set of types and this ensures the termination of drm .

Corollary 29. *For any type t and pattern p such that*

$$\forall \dot{x} \in \text{Acc}(p), \text{Op}(\dot{x}) \in \mathcal{O}_{\text{CDuce}}$$

the computation of $\Sigma; \square \vdash t // p$ terminates, where

$$\Sigma \equiv \{\dot{x} \mapsto \text{Init}(\dot{x}) \mid \dot{x} \in \text{Acc}(p)\}$$

Proof. We consider the type $t' = t \times t_0 \times \dots \times t_n$ where $t_i = \text{Init}(x_i)$. Thanks to Lemma 27, there is a plinth containing t' and therefore the rules in Figure 12 only generate a finite set of equations, the solution of which is a set of (mutually recursive) types.

Lemma 30. *All operators \mathcal{O}_{CDuce} have exact input.*

Proof. We consider each operator separately. The zipper projection operator Z_i is defined for every zipped type, so its accepted input is the type $(\mathbb{1})_{\top}$. Similarly, Z_i^L (resp., Z_i^R) is defined for every left (resp., right) projection zipper type, so its accepted input is the zipper type $L\mathbb{1} \cdot \top$ (resp., $R\mathbb{1} \cdot \top$). The projection operator π_i is defined for every pair, so its accepted input is the type $\mathbb{1} \times \mathbb{1}$. The top-level erasure rm is defined for every value and thus its accepted input is the type $\mathbb{1}$. The operator cons does not fail for any pair of values, so its accepted input is the type $\mathbb{1} \times \mathbb{1}$. Finally, snoc never fails for the first argument and fails only if its second argument is not a sequence. Therefore, the accepted input of snoc is the type $\mathbb{1} \times [\mathbb{1}^*]$.

Note that the set of operators we consider do not include, *e.g.*, the function application operator. Indeed, in general, this operator is not stable with respect to a given set of types. Of course, it does not mean that our calculus does not feature function application, but only —and this is a rather reasonable restriction— that function application cannot be used as an operator for accumulators.

C.1 Auxiliary Definitions and Lemmas

Before stating the soundness property for the whole language, we first state various auxiliary definitions and lemmas.

Lemma 31 (Strengthening). *Let Γ_1 and Γ_2 be two typing environments such that for any $x \in \text{dom}(\Gamma_1)$, we have $\Gamma_2(x) \leq \Gamma_1(x)$. If $\Gamma_1 \vdash e : t$, then $\Gamma_2 \vdash e : t$.*

Proof. By induction on the derivation of $\Gamma_1 \vdash e : t$. We simply introduce an instance of the subsumption rule below each instance of the [T-VAR] rule.

Lemma 32 (Admissibility of the intersection rule). *If $\Gamma \vdash e : t_1$ and $\Gamma \vdash e : t_2$, then $\Gamma \vdash e : t_1 \wedge t_2$.*

Proof. By induction on the structure of the two typing derivations.

Lemma 33. *Let Γ be a typing environment and e an expression that is well typed under Γ . Then the set*

$$S = \{t \in \mathcal{T} \mid \Gamma \vdash e : t \vee \Gamma \vdash e : \neg t\}$$

contains \emptyset and closed under \vee and \neg (and thus \wedge).

Proof. By definition, S is clearly closed under \neg . We have $\Gamma \vdash e : \mathbb{1} \simeq \neg \emptyset$ and thus $\emptyset \in S$. To show that S is closed under \vee , consider two types t_1 and t_2 in S . If $\Gamma \not\vdash e : t_1 \vee t_2$, then due to subsumption, we get $\Gamma \not\vdash e : t_1$ and $\Gamma \not\vdash e : t_2$. Because t_1 and t_2 are in S , we must have $\Gamma \vdash e : \neg t_1$ and $\Gamma \vdash e : \neg t_2$. By Lemma 32, we have $\Gamma \vdash e : \neg t_1 \wedge \neg t_2$ and $\neg t_1 \wedge \neg t_2 \simeq \neg(t_1 \vee t_2)$. Therefore, either $\Gamma \vdash e : t_1 \vee t_2$ or $\Gamma \vdash e : \neg(t_1 \vee t_2)$ holds, which completes the proof.

Lemma 34 (Substitution). *Let e, e_1, \dots, e_n be expressions, x_1, \dots, x_n distinct variables, t, t_1, \dots, t_n types, and Γ a typing environment. Then:*

$$\left\{ \begin{array}{l} \Gamma, (x_1 : t_1), \dots, (x_n : t_n) \vdash e : t \\ \forall i = 1..n, \Gamma \vdash e_i : t_i \end{array} \right\} \implies \Gamma \vdash e[e_1/x_1; \dots; e_n/x_n] : t$$

For simplicity, in this lemma, we do not distinguish variables from accumulators, writing both using the metavariable x .

Proof. By induction on the derivation of $\Gamma, (x_1 : t_1), \dots, (x_n : t_n) \vdash e : t$. We simply replace every instance of the rule [T-VAR] or [T-ACC] for variable x_i with a copy of the derivation of $\Gamma \vdash e_i : t_i$.

Definition 35. *We write $\llbracket t \rrbracket_{\mathcal{V}}$ for $\{v \mid \vdash v : t\}$.*

Lemma 36. *If $t \leq s$, then $\llbracket t \rrbracket_{\mathcal{V}} \subseteq \llbracket s \rrbracket_{\mathcal{V}}$. In particular, if $t \simeq s$, then $\llbracket t \rrbracket_{\mathcal{V}} = \llbracket s \rrbracket_{\mathcal{V}}$.*

Proof. Consequence of the subsumption rule.

Lemma 37. $\llbracket 0 \rrbracket_{\mathcal{V}} = \emptyset$.

Proof. We prove that $\vdash v : t$ implies $t \not\leq 0$ by induction on the typing derivation.

Lemma 38. $\llbracket t_1 \wedge t_2 \rrbracket_{\mathcal{V}} = \llbracket t_1 \rrbracket_{\mathcal{V}} \cap \llbracket t_2 \rrbracket_{\mathcal{V}}$.

Proof. By Lemma 36, $\llbracket t_1 \wedge t_2 \rrbracket_{\mathcal{V}} \subseteq \llbracket t_i \rrbracket_{\mathcal{V}}$ for $i = 1, 2$ and thus $\llbracket t_1 \wedge t_2 \rrbracket_{\mathcal{V}} \subseteq \llbracket t_1 \rrbracket_{\mathcal{V}} \cap \llbracket t_2 \rrbracket_{\mathcal{V}}$. Lemma 32 gives the opposite inclusion.

Lemma 39 (Inversion).

$$\begin{aligned} \llbracket t_1 \times t_2 \rrbracket_{\mathcal{V}} &= \{(v_1, v_2) \mid \vdash v_1 : t_1, \vdash v_2 : t_2\} \\ \llbracket b \rrbracket_{\mathcal{V}} &= \{c \mid b_c \leq b\} \\ \llbracket t \rightarrow s \rrbracket_{\mathcal{V}} &= \{\mu f^{(t_1 \rightarrow s_1; \dots; t_n \rightarrow s_n)}(x).e \in \mathcal{V} \mid \bigwedge_{i=1..n} t_i \rightarrow s_i \leq t \rightarrow s\} \\ \llbracket (t)_{\tau} \rrbracket_{\mathcal{V}} &= \{(w)_{\delta} \mid \vdash w : t, \vdash \delta : \tau\} \end{aligned}$$

Proof. For all four equalities, proving the \supseteq inclusion is straightforward. We prove the \subseteq inclusion by analyzing the typing derivation $\vdash v : t$, where t is instantiated to $t_1 \times t_2$, b , $t \rightarrow s$, or $(t)_{\tau}$ in each equality case.

Lemma 40. $\llbracket \neg t \rrbracket_{\mathcal{V}} = \mathcal{V} \setminus \llbracket t \rrbracket_{\mathcal{V}}$.

Proof. Note that $t \wedge \neg t \simeq 0$ and thus $\llbracket t \rrbracket_{\mathcal{V}} \cap \llbracket \neg t \rrbracket_{\mathcal{V}} = \llbracket t \wedge \neg t \rrbracket_{\mathcal{V}} = \llbracket 0 \rrbracket_{\mathcal{V}} = \emptyset$. Hence, it remains to prove that $\llbracket t \rrbracket_{\mathcal{V}} \cup \llbracket \neg t \rrbracket_{\mathcal{V}} = \mathcal{V}$, that is: $\forall v, \forall t, \vdash v : t \vee \vdash v : \neg t$. We prove this statement by induction over the pair (v, t) .

Lemma 41. $\llbracket t_1 \vee t_2 \rrbracket_{\mathcal{V}} = \llbracket t_1 \rrbracket_{\mathcal{V}} \cup \llbracket t_2 \rrbracket_{\mathcal{V}}$.

Proof. By Lemmas 36, 38, and 40.

Lemma 42. *Let p and φ respectively be a well-formed pattern and a zipper pattern. Let v and δ respectively be a closed value and a closed zipper. Then:*

$$\begin{aligned} (1) \quad & \forall \sigma, \text{Acc}(p) \subset \text{dom}(\sigma) \wedge \sigma; \delta^* \vdash v/p \not\rightsquigarrow \Omega \implies \vdash v : \lfloor p \rfloor \\ (2) \quad & \forall \sigma, \text{Acc}(\varphi) \subset \text{dom}(\sigma) \wedge \sigma \vdash \delta/\varphi \not\rightsquigarrow \Omega \implies \vdash \delta : \lfloor \varphi \rfloor \end{aligned}$$

Proof. By simultaneous induction on derivations. Note that values are inductively defined and patterns are regular and contractive.

Lemma 43. *Let p be a well-formed pattern, t a type such that $t \leq \lfloor p \rfloor$, x a variable such that $x \in \text{Var}(p)$, v a value, and $\delta^?$ a zipper and $\tau^?$ a zipper type such that $\vdash \delta^? : \tau^?$. Then for σ_0 such that $\text{Acc}(p) \subset \text{dom}(\sigma_0)$, the following hold:*

$$\begin{aligned} \exists v', (\vdash v' : t) \wedge (\sigma_0; \delta^? \vdash v'/p \rightsquigarrow \sigma, \gamma) \wedge (\gamma(x) = v) \\ \implies \exists t', (\tau^? \vdash t/p)(x) = t' \wedge \vdash v : t' \end{aligned}$$

Proof. By induction on a derivation of $\sigma_0; \delta^? \vdash v'/p \rightsquigarrow \sigma, \gamma$.

Lemma 44. *Let p be a well-formed pattern, t a type such that $t \leq \lfloor p \rfloor$, \dot{x} an accumulator such that $\dot{x} \in \text{Acc}(p)$, v a value, and $\delta^?$ a zipper and $\tau^?$ a zipper type such that $\vdash \delta^? : \tau^?$. Then for σ_0 and Σ_0 such that $\text{Acc}(p) \subset \text{dom}(\sigma_0) = \text{dom}(\Sigma_0)$ and $\forall \dot{x} \in \text{dom}(\sigma_0), \vdash \sigma_0(\dot{x}) : \Sigma_0(\dot{x})$, the following hold:*

$$\begin{aligned} \exists v', (\vdash v' : t) \wedge (\sigma_0; \delta^? \vdash v'/p \rightsquigarrow \sigma, \gamma) \wedge (\sigma(\dot{x}) = v) \\ \implies \exists \Sigma, \Sigma_0; \tau^? \vdash t//p = \Sigma \wedge \vdash v : \Sigma(\dot{x}) \end{aligned}$$

Proof. By induction on a derivation of $\sigma_0; \delta^? \vdash v'/p \rightsquigarrow \sigma, \gamma$.

C.2 Type Preservation

We now define what it means for an operator to be sound and then show that the set of operators we consider is sound.

Definition 45 (Sound operator). *An operator $(o, n_o, \overset{o}{\rightsquigarrow}, \overset{o}{\rightarrow})$ is sound if and only if for all $v_1, \dots, v_{n_o} \in \mathcal{V}$ such that $\vdash v_1 : t_1, \dots, \vdash v_{n_o} : t_{n_o}$, the following holds:*

$$\text{if } t_1, \dots, t_{n_o} \overset{o}{\rightarrow} s \text{ and } v_1, \dots, v_{n_o} \overset{o}{\rightsquigarrow} e, \text{ then } \vdash e : s$$

This allows us to finally state the soundness of the whole language (through type preservation, as usual).

Theorem 46 (Type preservation). *If all operators in the language are sound, then typing is preserved by reduction:*

$$\text{if } e \rightsquigarrow e' \text{ and } \vdash e : t, \text{ then } \vdash e' : t$$

In particular, $e' \neq \Omega$.

Proof. By induction on the derivation of $\vdash e : t$. We proceed by a case analysis on the last rule used in the derivation of $\vdash e : t$.

- **[T-CST]**: the expression e is a constant (value). It cannot be reduced, which contradicts the assumption, and thus the result follows.
- **[T-VAR]**: the expression e is a variable. It cannot be well typed under the empty context, which contradicts the assumption and thus the result follows.
- **[T-ACC]**: similar to the [T-VAR] case.
- **[T-ZIP-VAL]**: similar to the [T-CST] case.
- **[T-ZIP-EXPR]**: let $\vdash (e_0)_\bullet : (s)_\bullet$ where $e \equiv (e_0)_\bullet$, $t \equiv (s)_\bullet$, $\vdash e_0 : s$, and $s \leq \mathbb{1}_{\text{NZ}}$. Then, there exists an expression e'_0 such that $e_0 \rightsquigarrow e'_0$. We get $\vdash e'_0 : s$ by the induction hypothesis and then $\vdash (e'_0)_\bullet : (s)_\bullet$ by the rule [T-ZIP-EXPR].
- **[T-PAIR]**: let $\vdash (e_1, e_2) : t_1 \times t_2$ where $e \equiv (e_1, e_2)$, $t \equiv t_1 \times t_2$, $\vdash e_1 : t_1$, and $\vdash e_2 : t_2$. Then, there exists either an expression e'_1 such that $e_1 \rightsquigarrow e'_1$ or e'_2 such that $e_2 \rightsquigarrow e'_2$. If $e_1 \rightsquigarrow e'_1$, we get $\vdash e'_1 : t_1$ by the induction hypothesis and then $\vdash (e'_1, e_2) : t_1 \times t_2$ by the rule [T-PAIR]. The second case is similar to the first.
- **[T-SUB]**: there exists a type s such that $\vdash e : s$ and $s \leq t$. Since $e \rightsquigarrow e'$ by assumption, we have $\vdash e' : s$ by the induction hypothesis. Then, we get $\vdash e' : t$ by subsumption.
- **[T-OP]**: let $e \equiv o(e_1, \dots, e_{n_o})$ where $\forall i = 1..n_o$, $\vdash e_i : t_i$ and $t_1, \dots, t_{n_o} \xrightarrow{o} t$. There are two cases to consider:
 - (1) Suppose $o(e_1, \dots, e_i, \dots, e_{n_o}) \rightsquigarrow o(e_1, \dots, e'_i, \dots, e_{n_o})$ where $e_i \rightsquigarrow e'_i$. We get $\vdash e'_i : t_i$ by the induction hypothesis and then $\vdash o(e_1, \dots, e'_i, \dots, e_{n_o}) : t$ by the rule [T-OP].
 - (2) Suppose all e_i 's are values and $(e_1, \dots, e_{n_o}) \rightsquigarrow e'$. Since the operator $(o, n, \rightsquigarrow, \xrightarrow{o})$ is sound, we get $\vdash e' : t$ by Definition 45.
- **[T-FUN]**: similar to the [T-CST] case.
- **[T-MATCH]**: let $e \equiv \text{match } e_0 \text{ with } p_1 \rightarrow e_1 \mid p_2 \rightarrow e_2$. We have two cases to consider. First, suppose $e_0 \rightsquigarrow e'_0$. Then, we complete the proof by using the induction hypothesis and the rule [T-MATCH]. Now, suppose e_0 is a value. Then, either the matching e_0/p_1 succeeds or it fails and the matching e_0/p_2 succeeds. Here, we show only the proof of the first case; the second case is similar. From $\Gamma \vdash e : t$, we have the following assumptions:

- (1) $\vdash e_0 : s$
- (2) $s \leq \lambda p_1 \int \vee \lambda p_2 \int$
- (3) $s_1 \equiv s \wedge \lambda p_1 \int$
- (4) $s_2 \equiv s \wedge \neg \lambda p_1 \int$
- (5) $\Sigma_i \equiv \{\dot{x} \mapsto \text{Init}(\dot{x}) \mid \dot{x} \in \text{Acc}(p_i)\}$
- (6) $\Gamma_i \cup \Gamma'_i \vdash e_i : t_i$
- (7) $\Gamma_i \equiv \square \vdash s_i/p_i$
- (8) $\Gamma'_i \equiv \Sigma_i; \square \vdash s_i//p_i$
- (9) $t \equiv \bigvee_{\{i \mid s_i \neq 0\}} t_i$

- Let $\{\dot{x} \mapsto \text{Init}(\dot{x}) \mid \dot{x} \in \text{Acc}(p_1)\}; \square \vdash e_0/p_1 \rightsquigarrow \sigma, \gamma$.
- We get $\vdash e_0 : \lambda p_1 \int$ by Lemma 42.
- Then we get $\vdash e_0 : (s \wedge \lambda p_1 \int) \equiv s_1$ by Lemma 32.
- By rewriting (6), we have $\{(x : \Gamma_1(x)) \mid x \in \text{Var}(p_1)\} \cup \{(\dot{x} : \Gamma'_1(\dot{x})) \mid \dot{x} \in \text{Acc}(p_1)\} \vdash e_1 : t_1$.
- Furthermore, we have $\forall x \in \text{Var}(p_1)$, $\vdash \gamma(x) : \Gamma_1(x)$ by Lemma 43 and $\forall \dot{x} \in \text{Acc}(p_1)$, $\vdash \sigma(\dot{x}) : \Gamma'_1(\dot{x})$ by Lemma 44.

- Lemma 34 then gives us $\vdash e_1[\sigma; \gamma] : t_1$.
- Finally, by the rule [T-SUB], we get $\vdash e_1[\sigma; \gamma] : t$.

Since we made a runtime error (the special value Ω) explicit in the dynamic semantics, we do not need to show progress; showing type preservation is sufficient because Ω does not inhabit any type. All it remains to prove is that the operators we used are sound and whence deduce the soundness of the whole calculus.

Corollary 47. *The function application operator `app` is sound.*

Proof. Follows from Lemma 34.

Theorem 48. *The operators `app`, π_1 , π_2 , `drm`, `rm`, `cons`, and `snoc` are sound.*

Proof. Corollary 47 proves the soundness of the function application operator `app`. We prove the soundness of the operators π_1 , π_2 , `rm`, `cons`, and `snoc` by exploiting the fact that the types in the domains of their corresponding typing functions (especially for `cons` and `snoc`) are more precise than their exact input (defined in Lemma 30). We prove the soundness of `drm` in a similar way, using the soundness of `rm`.

Corollary 49. *Let \mathcal{O} be a set of operators $\{\text{app}, \pi_1, \pi_2, \text{drm}, \text{rm}, \text{cons}, \text{snoc}\}$. Our core calculus equipped with \mathcal{O} is sound in the sense that for any expression e , if $\vdash e : t$, then $e \not\rightarrow^* \Omega$.*

Proof. The result follows from Theorems 46 and 48.

D XPath Encoding

In this section, we give a translation of the so-called *navigational* XPath expressions into single patterns of the form $\text{axis}\{x \mid t\}$. We consider the following fragment of XPath.

Definition 50 (XPath query). *An XPath query is a finite term produced by the following grammar:*

$$\begin{aligned}
\text{path} &::= \text{step} \mid \text{path}/\text{step} \\
\text{step} &::= \text{axis}::\text{test}[\text{pred}] \\
\text{test} &::= l \mid * \\
\text{pred} &::= \text{path} \mid \text{not}(\text{pred}) \mid \text{pred} \text{ and } \text{pred} \mid \text{pred} \text{ or } \text{pred}
\end{aligned}$$

where axis produces the axes defined in Section 4, l ranges over XML element names (i.e., *atoms*), and $*$ is a wildcard test that denotes any label.

The semantics of XPath (as defined in [21]) relies on sets of nodes. Informally, given an initial set of nodes N_1 and a sequence of steps $s_1 / \dots / s_n$, the result of the evaluation of an XPath query is the set N_{n+1} obtained by

$$\begin{aligned}
s_1(N_1) &= N_2 \\
s_2(N_2) &= N_3 \\
&\vdots \\
s_n(N_n) &= N_{n+1}
\end{aligned}$$

An application of a step $a : l[p]$ to a set of nodes is computed as follows. First, for each node n in N , we compute the set N_a of nodes reachable through the axis a from n . Next, we filter N_a to keep only the set N_l of nodes whose label is l . Then, we keep only the set N_p of nodes of N_l for which the predicate p evaluates to **true**. Lastly, we *remove any duplicate* from N_p and return the nodes in document order. The truth value of predicates with respect to a node n is inductively defined as:

$$\frac{}{n \Vdash \text{path} : (\text{path}(\{n\}) \neq \emptyset)} \qquad \frac{n \Vdash p : b}{n \Vdash \text{not}(p) : \neg b}$$

$$\frac{n \Vdash p_1 : b_1 \quad n \Vdash p_2 : b_2}{n \Vdash p_1 \text{ or } p_2 : b_1 \vee b_2} \qquad \frac{n \Vdash p_1 : b_1 \quad n \Vdash p_2 : b_2}{n \Vdash p_1 \text{ and } p_2 : b_1 \wedge b_2}$$

Recalling the example given in the introduction (translated to the more concise syntax given above), the path $\text{desc}::\mathbf{a}[\text{not}(\text{anc}::\mathbf{b})]$ returns all descendants of the input that are labelled **a** and do *not* have an ancestor labelled **b**.

The biggest challenge in implementing the XPath semantics into patterns is to stick to the set-based semantics, without introducing internal node identifiers (that could be used to remove duplicates and sort the results in document order).⁵ This precludes giving a compositional, step-by-step semantics of XPath using the surface language. Indeed, consider the document:

$$d \equiv \langle \mathbf{a} \rangle [\langle \mathbf{a} \rangle [\langle \mathbf{a} \rangle []]]$$

Applying $\text{desc-or-self}::\mathbf{a}$ to d yields three intermediate results:

$$N = \{ \langle \mathbf{a} \rangle [\langle \mathbf{a} \rangle [\langle \mathbf{a} \rangle []]], \quad \langle \mathbf{a} \rangle [\langle \mathbf{a} \rangle []], \quad \langle \mathbf{a} \rangle [] \}$$

but applying $\text{desc-or-self}::\mathbf{a}$ again to N yields N itself, since the semantics is set-based.

The solution we propose is based on two key observations:

- i.* Given a predicate $[pred]$, we can write a CDuce type t such that for every value v , $v \Vdash pred : \text{true}$ if and only if $\vdash v : t$ (v has type t in CDuce);
- ii.* Any path $p \equiv s_1 / \dots / s_n$ can be put in the form of a predicate $[pred]$, such that the set of nodes selected by p is exactly the set of nodes for which $[pred]$ holds.

Using *(ii.)* we can put an XPath query in the form of a predicate and using *(i.)* we translate the predicate into a type t , and therefore express the whole XPath query as a pattern:

$$\text{desc-or-self}\{x \mid t\}$$

⁵ We could introduce a “set” type constructor exploiting internal node identifiers and new patterns for sets (besides sequences). However, this is merely the same as adding a separate layer for XPath on top of CDuce. Rather, we chose to encode XPath into CDuce patterns and thus benefit from the already existing static type system and efficient execution model of CDuce.

When this pattern is matched against a value v , it selects all the subtrees of v for which $[pred]$ holds, that is, all the subtrees returned by $p(\{v\})$.

Below, we first give a translation of XPath predicates into CDuce types and then a translation of XPath queries into predicates, both of which are well known in the literature (see Section 6 for more information).

Definition 51 (Encoding of XPath predicates into regular types). *Given a predicate ϕ (produced by the rule $pred$ of Definition 50), we define the mutually recursive functions T_{pred} , T_{path} , T_{step} , T_{axis} and T_{test} as follows:*

$$\begin{aligned}
T_{pred}(p) &= T_{path}(p) \\
T_{pred}(\phi_1 \text{ or } \phi_2) &= T_{pred}(\phi_1) \vee T_{pred}(\phi_2) \\
T_{pred}(\phi_1 \text{ and } \phi_2) &= T_{pred}(\phi_1) \wedge T_{pred}(\phi_2) \\
T_{pred}(\text{not } (\phi)) &= \neg T_{pred}(\phi) \\
\\
T_{path}(s) &= T_{step}(s, \mathbb{1}) \\
T_{path}(s/p) &= T_{step}(s, T_{path}(p)) \\
\\
T_{step}(a :: t[\phi], \tau) &= T_{axis}(a, \tau \wedge T_{test}(t)) \wedge T_{pred}(\phi) \\
\\
T_{test}(l) &= \langle l \rangle_{-} \\
T_{test}(\ast) &= \langle _ \rangle_{-} \\
\\
T_{axis}(\text{self}, \tau) &= \tau \\
T_{axis}(\text{child}, \tau) &= \langle _ \rangle [_ \ast \tau _ \ast] \\
T_{axis}(\text{desc-or-self}, \tau) &= \mu X. \tau \vee \langle _ \rangle [_ \ast X _ \ast] \\
T_{axis}(\text{desc}, \tau) &= T_{axis}(\text{child}, T_{axis}(\text{desc-or-self}, \tau)) \\
T_{axis}(\text{foll-sibling}, \tau) &= (_)_{L} (_ , [_ \ast \tau _ \ast]) \cdot \top \\
T_{axis}(\text{parent}, \tau) &= (_)_{L _ \cdot \mu X. ((R(\tau, _) \cdot (L _ \cdot \top \mid \bullet)) \mid R _ \cdot X)} \\
T_{axis}(\text{prec-sibling}, \tau) &= (_)_{L _ \cdot \mu X. ((R(\tau, _) \cdot R _ \cdot \top) \mid R _ \cdot X)} \\
T_{axis}(\text{anc}, \tau) &= (_)_{L _ \cdot \mu X. ((R(\tau, _) \cdot (L _ \cdot \top \mid \bullet)) \mid R _ \cdot X \mid L _ \cdot X)}
\end{aligned}$$

Example 52. The XPath predicate

$$p \equiv \text{parent} :: a \quad \text{or} \quad \text{desc-or-self} :: b / \text{child} :: c$$

is equivalent to the CDuce type t where:

$$\begin{aligned}
t &\equiv t_1 \vee t_2 \\
t_1 &\equiv (\mathbb{1})_{L _ \cdot \mu X. (R(\langle a \rangle \mathbb{1}) \cdot L _ \cdot \top \mid R _ \cdot X)} \\
t_2 &\equiv \mu X. (\langle b \rangle [\mathbb{1} \ast (\langle c \rangle \mathbb{1}) \mathbb{1} \ast] \vee \langle \mathbb{1}_{\text{basic}} \rangle [\mathbb{1} \ast X \mathbb{1} \ast])
\end{aligned}$$

Here, we see that the XPath operator **or** is translated into its set-theoretic counterpart (likewise for **and** and **not**). Upward path expressions are translated into zipper types, while downward path expressions are translated into recursive XML types. Furthermore, the chaining of steps is achieved by nesting the type obtained for the second step into the type obtained for the first step (see how $\langle c \rangle \mathbb{1}$ is embedded in t_2).

Now we formulate the translation of XPath queries of the form $s_1/\dots/s_n$ into an XPath predicate as the following lemma:

Lemma 53 (XPath to predicate translation [15]). *Let*

$$p \equiv a_1::l_1[p_1]/\dots/a_n::l_n[p_n]$$

be an XPath query. It can be translated into the following predicate:

$$p_n \text{ and self}::l_n/a_n^{-1}::l_{n-1}[p_{n-1}]/\dots/a_1^{-1}::*[\text{isroot}]$$

where a^{-1} is the inverse axis of a , defined as:

self^{-1}	= self	foll-sibling^{-1}	= prec-sibling
child^{-1}	= parent	desc-or-self^{-1}	= anc-or-self
parent^{-1}	= child	anc-or-self^{-1}	= desc-or-self
desc^{-1}	= anc	prec-sibling^{-1}	= foll-sibling
anc^{-1}	= desc		

and $\text{isroot} \equiv \text{not}(\text{parent}::)$.*

The key point of the translation is the fact that, for instance, the nodes selected by an XPath query $\text{child}::\mathbf{a}/\text{desc}::\mathbf{b}$ are all the nodes labelled \mathbf{b} that have an *ancestor* \mathbf{a} , the parent of which is the root of the document. In other words, the nodes selected by “ $\text{child}::\mathbf{a}/\text{desc}::\mathbf{b}$ ” are all the nodes for which the predicate “[$\text{self}::\mathbf{b}/\text{anc}::\mathbf{a}/\text{parent}::*[\text{isroot}]$]” holds.

To conclude the encoding of XPath, we add to our surface language one last extension to the syntax of expressions: “ e/path ” which is translated into:

$$\text{match } e \text{ with } \text{desc}\{x \mid t_{\text{path}}\} \rightarrow x \quad | _ \rightarrow []$$

where t_{path} is the type translation of the XPath query path . Note that the translation of e/path returns the results in document order since we use `snoc` for the accumulator used in the axis pattern $\text{desc}\{x \mid t_{\text{path}}\}$, thus faithfully implementing the semantics of XPath. With this extension, the CDuce version of the “*get_links*” function given in the introduction becomes as compact as in XQuery:

```
let get_links (page: <_>_) (print: <a>_ -> <a>_) : [ <a>_ * ] =
  transform page/desc::a[not(anc::b)] with x -> [ (print x) ]
```


E XQuery

It is well known (for example, see [24]) that full XPath expressions can be encoded using the XQuery fragment in Figure 7. In this section, we illustrate this with the following example:

Query in XQuery 3.0

```
1 for $i in
2   /descendant::time[ . = fn:dateTime()]/parent::elem
3   return $i
```

Query in XQ_H^+

```
4 for $d in document()
5   return
6     for $i in $d/descendant::time
7       return
8         switch
9           some $j in $i/self::node()
10            satisfies $j = fn:dateTime()
11            case true return for $k in $i/parent::elem
12              return $k
13          default return ()
```

The example query retrieves all `time` elements whose textual content is equal to the current time (returned by a built-in function) and return all the parent `elem` elements of such `time` elements. This query mixes within the path expression both navigational features (`descendant` and `parent` axes) with a data value test (it uses the string equality `=`). In contrast, the equivalent XQ_H^+ query makes an explicit, step by step navigation (navigational steps are highlighted in orange). Path predicates (between square brackets in XQuery) are translated into a conditional (the `switch_case` acting like an `if_then_else`) and use the Boolean construct `some $x in q_1 satisfies q_2` which evaluates to `true` if there exists an element of the sequence q_1 for which q_2 evaluates to `true` ($\$x$ is bound in turn to each element of q_1 and may appear free in q_2).

E.1 Typing Rules for XQ_H^+

$$\begin{array}{c}
\frac{\Gamma \vdash_{\text{xq}} q : s \quad s \leq t}{\Gamma \vdash_{\text{xq}} q : t} \text{typ-sub} \quad \frac{}{\Gamma \vdash_{\text{xq}} () : \text{'nil}} \text{typ-nil} \quad \frac{}{\Gamma \vdash_{\text{xq}} c : [c]} \text{typ-con} \\
\frac{\Gamma \vdash_{\text{xq}} q : t \quad t \leq [\mathbb{1}^*]}{\Gamma \vdash_{\text{xq}} \langle l \rangle q \langle /l \rangle : [\langle l \rangle t]} \text{typ-xml} \quad \frac{\Gamma \vdash_{\text{xq}} q_1 : [t_1^*] \quad \Gamma \vdash_{\text{xq}} q_2 : [t_2^*]}{\Gamma \vdash_{\text{xq}} q_1, q_2 : [t_1^* \ t_2^*]} \text{typ-seq} \\
\frac{x : t \in \Gamma}{\Gamma \vdash_{\text{xq}} \$x : t} \text{typ-var} \quad \frac{\begin{array}{c} \{y \mapsto \text{'nil}\}; \square \sim s // \text{axis}\{y \mid \mathbf{t}(\text{test})\} = \{y \mapsto t\} \\ \Gamma \vdash_{\text{xq}} x : [s^*] \quad t \leq [\mathbb{1}^*] \quad t' = \min\{t' \mid t \leq [t'^*]\} \end{array}}{\Gamma \vdash_{\text{xq}} \$x/\text{axis}::\text{test} : [t'^*]} \text{typ-path} \\
\frac{\Gamma \vdash_{\text{xq}} q_1 : [s^*] \quad \Gamma, x : [s] \vdash_{\text{xq}} q_2 : t \quad t \leq [\mathbb{1}^*]}{\Gamma \vdash_{\text{xq}} \text{for } \$x \text{ in } q_1 \text{ return } q_2 : t} \text{typ-for} \\
\frac{\Gamma \vdash_{\text{xq}} q : t \quad \begin{cases} t \not\leq \neg[c] \implies \Gamma \vdash_{\text{xq}} q_1 : s \\ t \leq [c] \implies \Gamma \vdash_{\text{xq}} q_2 : s \end{cases}}{\text{switch } q} \text{typ-case} \\
\frac{\Gamma \vdash_{\text{xq}} \text{case } c \text{ return } q_1 : s \quad \text{default return } q_2}{\Gamma \vdash_{\text{xq}} \text{case } c \text{ return } q_1 : s \quad \text{default return } q_2} \\
\frac{\Gamma \vdash_{\text{xq}} q : s \quad \begin{array}{c} t_1 = s \wedge \text{seq}(t) \quad \Gamma, x : t_1 \vdash_{\text{xq}} q_1 : t'_1 \\ t_2 = s \wedge \neg \text{seq}(t) \quad \Gamma \vdash_{\text{xq}} q_2 : t'_2 \end{array}}{\text{typeswitch } q} \text{typ-tcase} \\
\frac{\Gamma \vdash_{\text{xq}} \text{case } \text{tas } \$x \text{ return } q_1 : \bigvee_{\{i \mid t_i \neq 0\}} t'_i \quad \text{default return } q_2}{\Gamma \vdash_{\text{xq}} \text{case } \text{tas } \$x \text{ return } q_1 : \bigvee_{\{i \mid t_i \neq 0\}} t'_i \quad \text{default return } q_2} \\
\frac{\Gamma \vdash_{\text{xq}} q_1 : [s^*] \quad \Gamma, x : [s] \vdash_{\text{xq}} q_2 : [\text{bool}]}{\Gamma \vdash_{\text{xq}} \text{some } \$x \text{ in } q_1 \text{ satisfies } q_2 : [\text{bool}]} \text{typ-some} \\
\frac{\Gamma, x_1 : \text{seq}(t_1), \dots, x_n : \text{seq}(t_n) \vdash_{\text{xq}} q : \text{seq}(t)}{\Gamma \vdash_{\text{xq}} \text{fun } \$x_1 : t_1, \dots, \$x_n : t_n \text{ as } t. q : \text{seq}(t_1) \times \dots \times \text{seq}(t_n) \rightarrow \text{seq}(t)} \text{typ-fun} \\
\frac{\Gamma \vdash_{\text{xq}} q : t_1 \times \dots \times t_n \rightarrow t \quad \Gamma \vdash_{\text{xq}} q_i : t_i \quad (i = 1..n)}{\Gamma \vdash_{\text{xq}} q(q, \dots, q) : t} \text{typ-app}
\end{array}$$