

BiFluX: A Bidirectional Functional Update Language for XML

Hugo Pacheco

Cornell University, USA
hpacheco@cs.cornell.edu

Tao Zan

The Graduate University for Advanced
Studies, Japan
zantao@nii.ac.jp

Zhenjiang Hu

National Institute of Informatics, Japan
hu@nii.ac.jp

Abstract

Different XML formats are widely used for data exchange and processing, being often necessary to mutually convert between them. Standard XML transformation languages, like XSLT or XQuery, are unsatisfactory for this purpose since they require writing a separate transformation for each direction. Existing *bidirectional transformation* languages mean to cover this gap, by allowing programmers to write a single program that denotes both transformations. However, they often 1) induce a more cumbersome programming style than their traditionally unidirectional relatives, to establish the link between source and target formats, and 2) offer limited configurability, by making implicit assumptions about *how* modifications to both formats should be translated that may not be easy to predict.

This paper proposes a bidirectional XML update language called BiFLUX (BiDirectional Functional Updates for XML), inspired by the FLUX XML update language. Our language adopts a novel *bidirectional programming by update* paradigm, where a program succinctly and precisely describes *how* to update a source document with a target document, in an intuitive way, such that there is a unique “inverse” source query for each update program. BiFLUX extends FLUX with bidirectional actions that describe the connection between source and target formats. We introduce a core BiFLUX language, with a clear and well-behaved bidirectional semantics and a decidable static type system based on regular expression types.

Categories and Subject Descriptors D.3.1 [Programming Languages]: Formal Definitions and Theory; H.2.3 [Database Management]: Language—Data manipulation languages

Keywords XML, bidirectional programming, update languages

1. Introduction

Nowadays, various XML formats are widely used for data exchange and processing. Since data evolves naturally over time and is often replicated among different applications, it becomes frequently necessary to mutually convert between such formats. However, traditional XML transformation languages, like the XSLT and XQuery standards of the World Wide Web Consortium (W3C), are unsatisfactory for this purpose as they require writing a separate transformation for each direction.

Bidirectional transformation (BX) languages [11] mean to cover this gap, by allowing users to write a single program that can be executed both forwards and backwards, so that consistency between two formats can be maintained for free. A variety of bidirectional languages have emerged over the last 10 years to support bidirectional applications in the most diverse computer science disciplines [11], including functional programming, software engineering and databases. These languages come in different flavors, including many focused on the transformation of tree-structured data with a particular application to XML documents [7, 12, 14, 18, 21, 22, 25], and can be classified into three main paradigms. The first *relational* paradigm [7, 21] prescribes writing a declarative (non-deterministic) consistency relation between two formats, from which a suitable BX is automatically derived. The second *bidirectionalization* paradigm [12, 22, 24] asks users to write a transformation in a traditional unidirectional language, that plays the role of a functional consistency relation. The last *combinatorial* paradigm [14, 18, 25] encompasses the design of a domain-specific bidirectional language in which each combinator denotes a well-behaved BX, allowing users to write correct-by-construction programs by composition.

As most interesting examples of BXs are not bijective, there may be multiple ways to synchronize two documents into a consistent state, introducing ambiguity. Despite of this fact, bidirectional languages are typically designed to satisfy fundamental consistency principles, and support only a fixed set of synchronization strategies (out of a myriad possible) to translate a (non-deterministic) bidirectional specification—the syntactic description of a BX—into an executable BX procedure. This latent ambiguity often leads to unpredictable behavior, as users have limited power to configure and understand what a BX does from its specification. Even for combinatorial languages, that have the theoretical potential to fully specify the behavior of a BX [13, 27], their lower-level programming style requires significant effort and expertise from users to write intricate BXs via the composition of simple, concrete steps; they also scale badly for large formats, since one must explicitly describe how a BX transforms whole documents, including unrelated parts [25].

Intuitively, the goal of a BX is to translate updates on a target model into updates on a source model (and vice-versa) so that the updated models are kept consistent. As SQL stands for relational databases, a few high-level XML update languages [8, 15, 30] facilitate common modification operations over XML documents. Contrarily to XML transformation languages, XML update languages are well-suited for specifying small *in-place* changes to in a concise way, leaving all the remaining parts of a document unchanged.

In this paper, we propose a novel *bidirectional programming by update* paradigm, in which the programmer writes an update program that describes how to update a source model to embed information from a target model, and the system derives a query from source to target that evinces the consistency between both models. Such a *bidirectional update* allows to express the relationship between source and target models in a simple way—as in the

relational paradigm— by saying *which* related source parts are to be updated, but combined with additional actions that supply the missing pieces to tame the ambiguity in *how* target modifications are reflected—as in the combinatorial paradigm. For a wide class of BXs usually known as *lenses* [13], that have a data flow from *source* to *view*, this paradigm opens a new axis in the BX design space that enjoys a unique tradeoff between the declarative style of relational approaches and the stepwise style of combinatorial approaches. This paper demonstrates that a new family of *bidirectional update languages*, featuring an hybrid programming style, can render bidirectional programming more user friendly.

From a linguistic perspective, the main contribution of this paper is conceptual: we propose the idea of extending an update language with bidirectional features to write, directly and at a nice level of abstraction, a view update translation strategy which bundles all the necessary pieces to build a BX. Concretely, we design BiFLUX, a type-safe, declarative and expressive language for the bidirectional updating of XML documents that is deeply inspired by FLUX [8], a simple and well-designed functional XML update language. We lift unidirectional FLUX updates to bidirectional BiFLUX updates by imbuing them with an additional notion of view. Reading updates as BXs will motivate a few language extensions to original FLUX, and require a suitable bidirectional semantics and extra static conditions on BiFLUX programs to ensure that they build well-behaved BXs.

We demonstrate the usefulness of BiFLUX by illustrating typical examples of BXs written as bidirectional update programs. These help clarifying the stylistic differences to traditional bidirectional programming approaches, and substantiate that bidirectional update languages can combine a declarative language notation with a flexible and clear semantics. BiFLUX has been fully implemented and tested with many examples including those in this paper.

The rest of the paper is organized as follows. After briefly explaining the novel features of BiFLUX in Section 2, we show typical examples of BiFLUX programs in Section 3. Section 4 presents the core BiFLUX language that can be used to interpret high-level BiFLUX programs either as unidirectional or bidirectional updates, and Section 5 discusses the static typing and semantics of core BiFLUX. Section 6 formalizes the translation from high-level to core BiFLUX, Section 7 compares our approach with related work on bidirectional and XML programming, and Section 8 concludes with a synthesis of the main ideas and directions for future work.

2. A Bidirectional Update Language

Before giving concrete examples in Section 3, we start with a brief explanation of the features of BiFLUX and its informal semantics to show the big picture of our general framework.

2.1 BiFLUX syntax

We define the high-level syntax of BiFLUX in Figure 1 as a modest syntactic extension to FLUX; the new features are highlighted in green. Similarly to FLUX [8], BiFLUX enjoys a clear semantics and syntactic typechecking by translating programs into a canonical core language. Our examples assume an informal familiarity with commonplace XML technologies like XQuery expressions, XPath paths and XDuce-style regular expression types.

To have a taste of BiFLUX, imagine that we want to update the last author of a particular book with title ‘Querying XML’ in a database of books with type

```
books[book[title[string], author[string]+]+*
```

using a view of type *author*[string]. We can accomplish this by writing an update (with source *\$source* and view *\$view*):¹

```
UPDATE $source/books/book BY {
  REPLACE author[last()] WITH $view
} WHERE SOURCE title = "Querying XML"
```

FLUX-like syntax At first glance, bidirectional BiFLUX programs look just like regular FLUX programs. We omit the syntactic definitions of expressions *Expr*, paths *Path*, and patterns *Pat*. Variables *Var* are written *\$x*, *\$y*, etc. Statements *Stmt* include conditionals, composition, let-binding, case expressions or updates, which may be guarded by a **WHERE** clause that defines a set of conditions. Statements can be empty `{}` or parenthesized using braces `{Stmt}`. As in FLUX, *in-place* updates *Upd* can be *singular*, to update single trees, or *plural*, to update the children of each selected tree. Single insertions (**INSERT BEFORE/AFTER**) insert a value before or after each node selected by a path, while plural insertions (**INSERT AS FIRST/LAST INTO**) insert a value at the first or last position of the child-list of each selected node. Singular deletions (**DELETE**) delete each selected node, while plural deletions (**DELETE FROM**) delete their content. Single replacements (**REPLACE WITH**) replace each node selected by a path, while plural replacements (**REPLACE IN**) replace their content. Single updates **UPDATE BY** apply a statement to each tree in the result of a path.

Source and view matching The main difference in BiFLUX is that updates on sources carry an additional notion of view, what becomes syntactically evident with a new *non-in-place UPDATE FOR VIEW* operation that synchronizes a source sequence with a view sequence. Such synchronization can be configured by the programmer via a matching condition that aligns source and view nodes, and a triple of matching/unmatching clauses that describe the actions for individual source-view nodes. When two source and view nodes **MATCH**, a bidirectional statement is executed to update the source using the view; an unmatched view node (**UNMATCHV**) creates a new node in the source, either as a default or according to a unidirectional **CREATE** statement that provides a fresh source node to be normally updated with the existing view node; an unmatched source node (**UNMATCHS**) is **DELETED** by default, but we may **KEEP** it by providing a unidirectional statement describing how to invalidate the given **WHERE SOURCE** selection criteria. While **UPDATE FOR VIEW** statements are intrinsically bidirectional, the same BiFLUX syntax (e.g., **DELETE**) may be overloaded and denote either a bidirectional or a unidirectional update depending on the context. The rule is that all BiFLUX statements are bidirectional, except inside **UNMATCHS** or **UNMATCHV** clauses. An example that puts all these features to use is illustrated later in Figure 3.

Pattern matching Another significant difference to FLUX is the support for pattern matching. This is a very useful feature of XML transformation languages like XDuce [17] or CDuce [3], that allow matching tree patterns against the input data to transform it into an output of different shape. Typical XML update languages like XQuery! [15] or FLUX [8] do not support pattern matching, since it is not essential and may be more difficult to optimize, and they use solely paths to navigate to the portions of the input documents that are to be updated in-place. In BiFLUX, pattern matching can be used to guide the update based on the structure of the data (via **LET** and **CASE** statements), and is mostly useful for non-in-place updates that match source and view formats of different shapes.

Source/view/normal expressions Our language considers three kinds of source, view or normal variables. Expressions in **IF**, **LET**, **CASE** or **WHERE** clauses support additional tags to disambiguate if they refer to only the **SOURCE**, to only the **VIEW**, or to the global environment of an update. These tags can be ignored at first glance and will in fact be omitted in our examples, as they are inferable from context information for each update. Update procedures are omitted but can be easily added to the language.

¹ The position-dependent *last()* XPath function is actually not supported in our path expressions, and is desugared in BiFLUX using pattern matching.

$ \begin{aligned} \text{Stmt} & ::= \text{Upd} [\text{WHERE } \text{Conds}] \mid \text{Stmt} ; \text{Stmt} \mid \{ \text{Stmt} \} \mid \{ \} \\ & \mid \text{IF } \text{Tag } \text{Expr} \text{ THEN } \text{Stmt} \text{ ELSE } \text{Stmt} \\ & \mid \text{LET } \text{Tag } \text{Pat} = \text{Expr} \text{ IN } \text{Stmt} \\ & \mid \text{CASE } \text{Tag } \text{Expr} \text{ OF } \{ \text{Cases} \} \\ \text{Upd} & ::= \text{INSERT (BEFORE} \mid \text{AFTER)} \text{ PatPath VALUE Expr} \\ & \mid \text{INSERT AS (FIRST} \mid \text{LAST)} \text{ INTO PatPath VALUE Expr} \\ & \mid \text{DELETE [FROM] PatPath} \mid \text{REPLACE [IN] PatPath WITH Expr} \\ & \mid \text{UPDATE PatPath BY Stmt} \\ & \mid \text{UPDATE PatPath BY VStmt FOR VIEW PatPath [Match]} \\ & \mid \text{KEEP PatPath} \mid \text{CREATE VALUE Expr} \\ \text{Conds} & ::= \text{Tag Expr} \text{ [; Conds]} \mid \text{Tag Var} \text{ := Expr} \text{ [; Conds]} \end{aligned} $	$ \begin{aligned} \text{Cases} & ::= \text{Pat} \rightarrow \text{Stmt} \mid \text{Cases} \text{ '}' \text{ Cases} \\ \text{VStmt} & ::= \{ \text{VStmt} \} \mid \text{VUpd} \\ & \mid \text{VUpd} \text{ '}' \text{ VUpd} \\ \text{VUpd} & ::= \text{MATCH} \rightarrow \text{Stmt} \\ & \mid \text{UNMATCHS} \rightarrow \text{Stmt} \\ & \mid \text{UNMATCHV} \rightarrow \text{Stmt} \\ \text{Match} & ::= \text{MATCHING BY Path} \\ & \mid \text{MATCHING SOURCE BY Path} \\ & \qquad \qquad \text{VIEW BY Path} \\ \text{PatPath} & ::= [\text{Pat IN}] \text{ Path} \\ \text{Tag} & ::= [\text{SOURCE} \mid \text{VIEW}] \end{aligned} $
--	--

Figure 1: Concrete syntax of BiFLUX updates.

2.2 Informal semantics and general framework

In general, a BiFLUX update is executed for a particular source and view as follows: a source path is evaluated over the current source, yielding a *source focus selection*, to be recursively updated using a *view focus selection* computed by evaluating a view path over the current view, until all the view information is embedded into the source. View and source focus selections denote the mutable parts of the source and view trees that can be updated and used by the update, and subexpressions of the update may restrict the focuses.

Despite the emphasis is on writing updates, BiFLUX programs have a bidirectional interpretation. They can be read as 1) an *update function* $U(s, v') = s'$ that updates a source s into a new source s' which contains a given view v' , or 2) a *query function* $Q(s) = v$ that computes a view v from a given source s ; these functions may be partial². For the example at the beginning of this section, (assuming that books are uniquely identified by their titles) the corresponding query function is semantically equivalent to the XPath expression that returns the last author of the respective source book:

```
$source/books/book[title="Querying XML"]/author[last()]
```

Our language is carefully designed to ensure that the inferred relationship between sources and views is deterministic, so that capturing it by a query function is appropriate. In other words, there exists a unique query function for each update program written in our language. Moreover, its bidirectional semantics satisfies two basilar synchronization properties; that an update U consistently embeds view information to the source, without view side-effects:

$$U(s, v') = s' \Rightarrow Q(s') = v' \quad \text{UPDATEQUERY}$$

and that it does not update already consistent sources:

$$Q(s) = v \Rightarrow U(s, v) = s \quad \text{QUERYUPDATE}$$

These two properties are commonly known as the well-behavedness laws of lenses in the bidirectional programming community [11].

For an example of a FLUX update that is not (syntactically) valid as a BiFLUX update, imagine that we had written instead:

```
UPDATE $source/books/book BY {
  INSERT AS LAST INTO author VALUE $view
} WHERE SOURCE title = "Querying XML"
```

This update function is not idempotent on sources, since it always inserts the view as an extra author of the source book, violating QUERYUPDATE: even when the view is already an author of the source book, a new duplicated author is inserted. Such class of valid BiFLUX updates can be statically checked, namely as the programs for which update normalization and typechecking succeed.

The general architecture of our bidirectional updating framework is illustrated in Figure 2. A BiFLUX program is evaluated in two

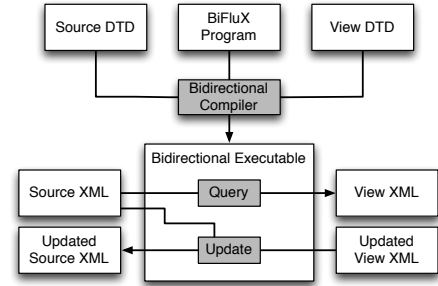


Figure 2: Architecture of the BiFLUX framework.

stages. First, it is statically compiled against a source and a view schema (represented as DTDs), producing a bidirectional executable. The generated executable can then be evaluated bidirectionally for particular XML documents conforming to the DTDs: in forward mode as a query Q , or in backward mode as an update U .

3. Examples

This section illustrates how writing a bidirectional update feels like to a programmer, through a series of BX examples using BiFLUX:

Example 3.1 demonstrates how the new bidirectional notation, in combination with ordinary unidirectional FLUX updates, can be used to intuitively describe a bidirectional update;

Example 3.2 illustrates nested bidirectional updates and flexible non-local matching behavior;

Example 3.3 showcases pattern matching and recursive procedures.

3.1 Institutional address book example

Consider a typical address book represented by the DTD from Figure 4. An address book contains a list of persons, each possessing a name, a list of emails and an optional telephone number. Let us start with the following XML address book with three persons:

```
<addrbook><person><name>Hugo Pacheco</name>
  <email>hpacheco@nii.ac.jp</email>
  <email>hpacheco@gmail.com</email></person>
  <person><name>John Doe</name>
  <email>doe@domain.com</email></person>
  <person><name>Zhenjiang Hu</name><email>zh@nii.ac.jp
  </email><tel>+81-3-4212-2530</tel></person>
</addrbook>
```

On the other hand, the NII's administrative services may keep only a view with the name and institutional address of employees (people with an email at "nii.ac.jp"), as shown in the DTD from Figure 5.

² We represent the undefined value as \perp , failure of a partial function f as $f v = \perp$ and partial inclusion ($v \sqsubseteq v'$) as $v \neq \perp \Rightarrow v = v'$.

```

PROCEDURE niibook(SOURCE $source AS s:addrbook, VIEW $view AS v:niibook) =
UPDATE $source/addrbook/person BY {
MATCH -> REPLACE email[ends-with(text(), 'nii.ac.jp')] [1] WITH $email'
| UNMATCHV -> CREATE VALUE <person> <name/> <tel>+81-3-4212-2000</tel> </person>
| UNMATCHS -> KEEP . ; DELETE email[ends-with(text(), 'nii.ac.jp')]
} FOR VIEW employee[$name AS v:name, $email' AS v:email] IN $view/niibook/employee
MATCHING BY name WHERE email[ends-with(text(), 'nii.ac.jp')]

```

Figure 3: BiFLUX update for the institutional address book example.

We can easily write a bidirectional update in BiFLUX to synchronize these two formats, as illustrated in Figure 3³. The root procedure *niibook* takes as arguments the root source and view variables. It focuses on a sequence of source persons by traversing down the path *\$source/addrbook/person*, selecting only those that have NII emails, and focuses on a sequence of view employees by traversing up the (injective) path *\$view/niibook/employee*. Elements in the two sequences are matched by their names. For matching person-employee pairs, the person's first NII email is updated with the employee's email. If a new unmatched employee exists in the view, a new person with a default telephone is created in the source (inheriting its view name and email). If an old unmatched person exists in the source, all its NII emails are deleted.

The unique query for this example simply keeps the first institutional email of each person working at the NII:

```

<niibook><employee><name>Hugo Pacheco</name>
  <email>hpacheco@nii.ac.jp</email></employee>
<employee><name>Zhenjiang Hu</name>
  <email>zh@nii.ac.jp</email></employee>
</niibook>

```

The update function is more interesting. For instance, if we add Tao (in alphabetical order) as a new NII employee, fix Zhenjiang's email and delete Hugo, we get an updated source where John is left unchanged, Tao is created with a default telephone number (as his name does not match any name in the original source), Zhenjiang's NII email is updated and the NII email of Hugo is deleted:

```

<addrbook><person><name>Hugo Pacheco</name>
  <email>hpacheco@gmail.com</email></person>
<person><name>John Doe</name>
  <email>doe@domain.com</email></person>
<person><name>Tao Zan</name><email>zantao@nii.ac.jp
  </email><tel>+81-3-4212-2000</tel></person>
<person><name>Zhenjiang Hu</name><email>hu@nii.ac.jp
  </email><tel>+81-3-4212-2530</tel></person>
</addrbook>

```

Note that if we queried the updated address book again, we would get a view with the names and NII emails of only Tao and Zhenjiang.

As a BX side note, this precise behavior can not be achieved using the existing typical declarative BX languages, which are not designed with fine user control in mind; alignment-aware combinatorial BX languages could be tailored to produce similar behavior, but getting it right requires much higher effort and expertise.

3.2 Social address book example

In response to the Web 2.0 movement, consider that our address book now supports groups of people according to their social relationships. The new DTD corresponds to the following regular expression type:

```
addrbook[group[name[string], person[...]]*]*
```

Consider our address book with people now classified into groups:

```
<addrbook><group><name>coworkers</name>
```

³The names *s:elem* and *v:elem* are BiFLUX type variables that refer to the types of source and view elements declared in the respective DTDs.

Figure 4: A simple address book DTD.

```

<!DOCTYPE niibook [ <!ELEMENT niibook (employee)*>
<!ELEMENT employee (name,email)>
<!ELEMENT name (#PCDATA)> <!ELEMENT email (#PCDATA)> ]>

```

Figure 5: A NII address book DTD.

```

<person><name>Hugo Pacheco</name>...</person>
<person><name>Zhenjiang Hu</name>...</person></group>
<group><name>friends</name><person>
  <name>John Doe</name>...</person></group>
</addrbook>

```

This time, a social media application may only be interested in the groups and names of people, according to a view schema:

```
socialbook[group[name[string], person[name[string]]*]*]
```

A bidirectional update that synchronizes address and social books is written in Figure 6. It starts by matching groups, proceeding recursively for persons within groups (with default behavior for unmatched groups). Inside, for source-view persons matching on their name, no update is necessary. For each unmatched view person, we attempt to retrieve its address book information from any other group in the original source, or otherwise create a default person.

The corresponding query function produces a view with the same structure but showing only names of groups and persons. For the update function, imagine that we modify the view by reordering the two groups, changing Hugo's group and creating a new group for family members:

```

<addrbook><group><name>friends</name>
  <person><name>Hugo Pacheco</name></person></group>
<group><name>coworkers</name><person>
  <name>Zhenjiang Hu</name></person></group>
<group><name>family</name></group>
</addrbook>

```

The correspondingly updated source is as follows:

```

<addrbook><group><name>friends</name>
  <person><name>Hugo Pacheco</name>...</person>
  <person><name>John Doe</name>...</person></group>
<group><name>coworkers</name><person>
  <name>Zhenjiang Hu</name>...</person></group>
<group><name>family</name></group>
</addrbook>

```

Since Hugo changed from group coworkers to friends, he is considered an unmatched view person under his new group. Our bidirectional update avoids his original address details to be lost, by looking them up in all groups, instead of only in Hugo's original group (what would be the default behavior). An analogous example motivates an extension to the alignment-aware language of [1].

3.3 Bookmark example

For a different bidirectional updating example, consider the conversion between two popular browser bookmark formats studied in [21]. Netscape stores its bookmarks in an HTML format (Figure 7), while


```

PROCEDURE socialbook(SOURCE $source AS s:addrbook, VIEW $view v:socialbook)
= UPDATE $source/addrbook/group BY { MATCH ->
  UPDATE $person IN $persons BY { MATCH -> {}
  | UNMATCHV -> LET $old = $source/addrbook/group/person IN
    LET $oldperson = $old[name/text() = $person'/name/text()][1] IN
    IF $oldperson THEN CREATE VALUE $oldperson ELSE {}
  } FOR VIEW $person' IN $persons' MATCHING BY name
} FOR VIEW group[$name AS v:name, $persons AS v:person*]
  IN $view/socialbook/group MATCHING BY name
<!DOCTYPE html[<!ELEMENT html(head,body)>
<!ELEMENT head (#PCDATA)>
<!ELEMENT body (h1,dl)>
<!ELEMENT h1 (#PCDATA)>
<!ELEMENT dl ((dt|dd)*)>
<!ELEMENT dt (a)>
<!ELEMENT a (href,#PCDATA)>
<!ELEMENT dd (h3,dl)>
<!ELEMENT h3 (#PCDATA)> ]>

```

Figure 6: BiFLUX update for the social address book example.

Figure 7: A Netscape bookmark format DTD.

```

PROCEDURE top(SOURCE $html AS s:html,VIEW $xbel AS v:xbel) =
UPDATE html[head[String], body[$h1 AS s:h1, dl[$nc AS (s:dt|s:dd)*]] IN $html BY
{ REPLACE IN $h1 WITH $t ; contents($nc,$xc) }
FOR VIEW xbel[title[$t AS String], $xc AS (v:bookmark|v:folder)*] IN $xbel

PROCEDURE contents(SOURCE $nc AS (s:dt|s:dd)*,VIEW $xc AS (v:bookmark|v:folder)* = UPDATE $nc BY { CASE $v OF {
  bookmark[href[$url AS String], title[$title AS String]] -> REPLACE . WITH <dt><a><href>{$url}</href>{$title}</a></dt>
  | folder[title[$title AS String], $fxc AS (v:bookmark|v:folder)*] -> REPLACE IN h3 WITH $title ; contents(dl/*,$fxc)
} } FOR VIEW $v IN $xc

```

Figure 8: BiFLUX update for the bookmark example.

```

<!DOCTYPE xbel[<!ELEMENT xbel(title,(bookmark|folder)*)>
<!ELEMENT title(#PCDATA)><!ELEMENT bookmark(href,title)>
<!ELEMENT folder(title,(bookmark|folder)*)> ]>

```

Figure 9: A XBEL bookmark format DTD.

the XBEL open XML bookmark exchange format opts for a loosely equivalent representation (Figure 9). Both formats contain a general title (h1 or title) and a sequence of bookmarks (dt or bookmark) or folders (dd or folder), where folders may recursively contain sequences of bookmarks or folders.

The original biXid transformation [21] relies on pattern matching to decompose the source and target formats and can be replicated in BiFLUX as shown in Figure 8. The *top* procedure decomposes the source into head and body (with a title *\$h1* and a sequence *\$nc* of dts or dds) and the view into a title *\$t* and a sequence *\$xc* of bookmarks or folders. Then *top* replaces the source *\$h1* with *\$t* and invokes *contents* to update the remaining sequences. The *contents* procedure makes use of a case expression to match source and view bookmarks and folders: for a view *bookmark*, we generate a source *dt* element with the bookmark’s *href* and *title*; for a view *folder*, we generate a source *dd* element with the folder’s *title* and a *dl* with recursively computed contents.

Finally, we can run our BiFLUX program as a query that converts from Netscape to XBEL, or as an update that converts from XBEL to Netscape. This is not the best showcase example of BiFLUX, since the source and view bookmarks are almost in bijective correspondence and there is small ambiguity to mitigate in the update. Nonetheless, note that our BiFLUX program will preserve the original Netscape header, while the analogous biXid program would simply generate default data for unrelated parts.

4. Core Language

The high-level language presented in the previous sections follows a verbose and natural syntax that is convenient for users, but its operations are overlapping, complex and hard to typecheck. Following the design of FLUX and as standard for many other languages, we

introduce a core update language of canonical operations whose semantics and typing rules are easier to define and manipulate.

4.1 XML values and regular expression types

As several other XML processing languages [8, 9, 17], we consider a type system of regular expression types with structural subtyping⁴:

Atomic types $\alpha ::= \text{bool} \parallel \text{string} \parallel n[\tau]$
Sequence types $\tau ::= \alpha \parallel () \parallel \tau \mid \tau' \parallel \tau, \tau' \parallel \tau^* \parallel X$

Atomic types $\alpha \in \text{Atom}$ are primitive booleans, strings or labeled sequences $n[\tau]$. *Sequence types* $\tau \in \text{Type}$ are defined using regular expressions, including empty sequence $()$, alternative choice $\tau \mid \tau'$, sequential composition τ, τ' , iteration τ^* or type variables X ; choice and composition are right-nested. We define the usual $\tau^+ = \tau, \tau^*$ and $\tau^? = \tau \mid ()$. Types can also be recursively defined:

Type definitions $\tau_D ::= \alpha \parallel () \parallel \tau_D \mid \tau'_D \parallel \tau_D, \tau'_D \parallel \tau_D^*$
Type signatures $E ::= \cdot \parallel E, \text{type } X = \tau_D$

Type definitions τ_D are sequences with no top-level variables (to avoid non-label-guarded recursion [9]). A *type signature* E is a set of named type definitions of the form $X = \tau_D$, and is well-formed if no two types have the same name and all type variables in definitions are declared in E . We write $E(X)$ for the type bound to X in E . Hereafter, we will assume the signature E to be fixed.

In traditional XML-centric approaches [9, 17], values are encoded using a uniform representation that does not record the structure that types impose on values. This “flat” representation is economical and simplifies subtyping, but makes it harder to realize that a value belongs to a type and therefore to integrate regular expression features into functional languages with non-structural type equivalence, such as Haskell or ML. In this paper, we instead consider a structured representation of values (in line with values of algebraic data types) that keep explicit annotations which, in a way, witness how to parse a flat value as an instance of a type [23]:

Atomic values $t ::= \text{true} \mid \text{false} \mid w \mid n[v]$
Forest values $v ::= t \mid () \mid L v \mid R v \mid (v, v) \mid [v_0, \dots, v_n]$

⁴ We use \parallel for syntax alternatives in the type grammar to prevent confusion.

Atomic values $t \in Tree$ can be `true`, `false` $\in Bool$, strings $w \in \Sigma^*$ (for some alphabet Σ), or singleton trees $n[v]$ with a node label n . Forest values $v \in Val$ include the empty sequence $()$, left- L v or right- R v tagged choices, binary sequences (v, v) and lists of arbitrary length $[v_0, \dots, v_n]$. The semantics of a type τ denotes a set of values $\llbracket \tau \rrbracket$ that is defined as the minimal solution (formally the least fixed point [17]) of the following set of equations:

$$\begin{aligned} \llbracket \text{bool} \rrbracket &\triangleq \{\text{true}, \text{false}\} & \llbracket n[\tau] \rrbracket &\triangleq \{n[v] \mid v \in \llbracket \tau \rrbracket\} \\ \llbracket \text{string} \rrbracket &\triangleq \Sigma^* & \llbracket X \rrbracket &\triangleq \llbracket E(X) \rrbracket \\ \llbracket \tau, \tau' \rrbracket &\triangleq \{(v, v') \mid v \in \llbracket \tau \rrbracket, v' \in \llbracket \tau' \rrbracket\} & \llbracket () \rrbracket &\triangleq \{()\} \\ \llbracket \tau \mid \tau' \rrbracket &\triangleq \{L v \mid v \in \llbracket \tau \rrbracket\} \cup \{R v \mid v \in \llbracket \tau' \rrbracket\} \\ \llbracket \tau^* \rrbracket &\triangleq \{[v_0, \dots, v_n] \mid v_0, \dots, v_n \in \llbracket \tau \rrbracket, n \geq 0\} \end{aligned}$$

In our context, values in the type semantics preserve the type structure. We will denote flat values $ft \in FTree$ and $fv \in FVal$ (dropping left/right tags, parenthesis and list brackets) by:

$$\begin{aligned} \text{Flat atomic values} \quad ft &::= \text{true} \mid \text{false} \mid w \mid n[fv] \\ \text{Flat forest values} \quad fv &::= () \mid ft, fv \end{aligned}$$

and introduce a function $flat : Val \rightarrow FVal$ that ignores markup. We denote the usual flat semantics of a type τ as $\llbracket \tau \rrbracket_{flat}$.

4.2 Core expression, path and pattern language

In BiFLUX, updates instrumentally use XQuery expressions, XPath paths and XDuce patterns to manipulate XML data. This subsection succinctly describes their syntax in our core language, but is not essential for our design and may be skipped on a first reading.

We write *expressions* e in a minimal XQuery-like language that is a variant of the μXQ core language proposed in [9]:

$$\begin{aligned} e &::= () \mid e, e' \mid n[e] \mid \text{let } pat = e \text{ in } e' \mid p \mid e \approx e' \\ &\mid \text{if } e \text{ then } e' \text{ else } e'' \mid \text{for } x \text{ in } e \text{ return } e' \end{aligned}$$

Despite expressions can be used in updates rather indiscriminately, in BiFLUX only a particular subset of the expression language is suitable for denoting foci. Therefore, we differentiate *paths* p in a core path language that represents a minimal dialect of XPath:

$$\begin{aligned} p &::= \text{self} \mid \text{child} \mid \text{dos} \mid :: nt \mid \text{where } e \mid p / p' \\ &\mid x \mid w \mid \text{true} \mid \text{false} \mid F(\vec{e}) \\ nt &::= n \mid \text{text}() \mid \text{node}() \end{aligned}$$

To simplify the formal treatment, we consider nodetests $::nt$ that apply to atomic values and where clauses $\text{where } e$ that filter values satisfying an expression e . We write the syntactic sugar $p :: nt \triangleq p / ::nt$ and $p[e] \triangleq p / \text{where } e$. XPath-like traversals can be defined as $. = \text{self} :: \text{node}()$, $p / n \triangleq p / \text{child} :: a$ and $p // n = p / \text{dos} :: \text{node}() / \text{child} :: a$.

In contrast to μXQ , our core expression language supports pattern expressions in let bindings: $\text{let } pat = e \text{ in } e'$ matches a pattern pat against the result of an expression e and then executes an expression e' that may refer to the variables bound by pat . We consider a language of XDuce-style *patterns* pat [17]:

$$pat ::= x \text{ as } \tau \mid \tau \mid () \mid n[pat] \mid pat, pat'$$

Note that we impose a simple but strong syntactic *linearity* restriction on patterns (no alternative choice, no Kleene star) to ensure that matching a value against a pattern binds each variable exactly once. Less severe linearity restrictions are actually known [16], but these simple patterns suffice for our practical needs. Also worth noting is that we require every variable to be annotated with a type. This simplifies our design, but will in turn increase the number of (often unnecessary) annotations in our bidirectional update programs. We see it as an orthogonal problem that can be mitigated using existing tree-based type inference algorithms [32].

4.3 Tripartite environments

A *type environment* Γ consists of a set of bindings $x : \tau$ of variables to types. An *environment* is a function $\gamma : Var \rightarrow Val$ from variables to values. As usual, we assume variable names in an environment to be distinct. We write $\Gamma(x)$ and $\gamma(x)$ for the type and value of a variable; $\Gamma[x : \tau]$ and $\gamma[x := v]$ add a new variable to an environment. We say that γ has type Γ , written $\gamma : \Gamma$, if $\gamma(x) : \Gamma(x)$ for all $x \in dom(\Gamma)$. Fresh environments are denoted by the empty set \emptyset .

Since our language is bidirectional, we will consider three kinds of variables (and environments): *source variables*, that are accessible from the current source; *view variables*, that are accessible from the current view; and *normal variables*, that are accessible from the global environment of the update and independent from the source or the view. We will talk about source/view/normal paths or expressions, that may only refer to source variables, view variables or any variable, respectively; *non-source* expressions may refer to view and normal variables simultaneously.

To describe source or view environments, we introduce *record types* ν that denote sequences of variable-annotated types:

$$\nu ::= x : \tau, \nu \mid ()$$

As an abuse of notation, we see a record type ν as an ordinary type (by forgetting variable names) or as a type environment; we cast a conforming value v into an environment $\gamma_{v:\nu}$.

4.4 Core update language

Unlike conventional XML update languages, our core update language comprises two kinds of update statements: unidirectional FLUX updates, interpreted as *arrows* [20] that modify a document in-place, changing its schema; and bidirectional BiFLUX updates, interpreted as BXs [27] that update a source document given a view document or query a source document to compute its view fragment, for fixed source and view schemas. Our core *unidirectional updates* u are adapted from the core FLUX update language [8]:

$$\begin{aligned} u &::= \text{skip} \mid u; u' \mid \text{insert } e \mid \text{delete} \\ &\mid \text{if } e \text{ then } u \text{ else } u' \mid \text{case } e \text{ of } \vec{pat} \rightarrow \vec{u} \\ &\mid p[u] \mid \text{left}[u] \mid \text{right}[u] \mid \text{children}[u] \end{aligned}$$

These include standard operations such as the no-op `skip`, sequential composition, conditionals or case expressions. The basic operations are `insert` e , that inserts a value given an empty sequence as focus; and `delete`, that replaces any value with the empty sequence. We can also apply an update in a specific *direction* (that traverses down a path p , moves to the left or right of a value, or focuses on the children of a labeled node).

Core *bidirectional updates* b are denoted by the grammar:

$$\begin{aligned} b &::= \text{skip} \mid \text{fail} \mid b_1; b_2 \mid \text{view } x := e \text{ in } b \\ &\mid P(\vec{p}_s, \vec{e}_v, \vec{e}) \mid p[b] \mid [b]e \mid \text{replace} \mid \text{iter } b \\ &\mid \text{ifS } e \text{ then } b \text{ else } b' \mid \text{caseS } p \text{ of } \vec{pat} \rightarrow \vec{b} \\ &\mid \text{ifV } e \text{ then } b \text{ else } b' \mid \text{caseV } e \text{ of } \vec{pat} \rightarrow \vec{b} \\ &\mid \text{if } e \text{ then } b \text{ else } b' \mid \text{case } e \text{ of } \vec{pat} \rightarrow \vec{b} \\ &\mid \text{alignpos } e_s \ b \ c \ r \mid \text{alignkey } e_s \ p_s \ p_v \ b \ c \ r \end{aligned}$$

Here, P is the name of a BiFLUX *procedure*. A procedure is defined as a declaration $P(\vec{x} : \vec{\tau}) : \nu_s \Leftrightarrow \nu_v \triangleq s$, meaning that P takes a vector of parameters \vec{x} of types $\vec{\tau}$ and builds a BX between a source of type ν_s and a view of type ν_v . Accordingly, a procedure call $P(\vec{p}_s, \vec{e}_v, \vec{e})$ takes three kinds of arguments: source paths \vec{p}_s ; non-source expressions \vec{e}_v ; and normal expressions \vec{e} . We collect procedure declarations into a set Δ , that we will assume to be fixed. Procedures may also be recursive.

The bidirectional operation `skip` keeps the source unchanged for an empty view; for a non-empty view, we must `fail` to update the source (as UPDATEQUERY precludes that all view information must

be used to update the source). Composition $b_1; b_2$ updates the source with b_1 , and then runs b_2 over the updated source. The statement `view $x := e$ in b` models a view dependency, by stating that a view variable x can be computed using the view expression e (written $\$x := e$ in the high-level BIFLUX language), and then runs b using the remaining view. Bidirectional statements may also change the current source or view, by focusing down a source path and updating the resultant source ($p[b]$), or by evaluating a non-source expression “backwards” from the view and updating the source with the resultant view ($[b]e$). The basic `replace` operation embeds the view into the source, while `iter b` embeds the same view into each tree in a source forest. As before, we consider three kinds of bidirectional conditionals and case expressions, whose expressions or paths are respectively source, view/non-source or normal variables.

The two special alignment statements update a source sequence using a view sequence. They receive a filtering source expression e_s , and match source elements satisfying e_s with view elements by position (`alignpos`) or by keys (`alignkey`), defined by two source (p_s) and view paths (p_v). Then, a bidirectional statement b processes matched source-view elements, a create statement c instructs how to create a suitable source to match with an unmatched view, and a recover statement r denotes how to process an unmatched source. Create statements c are simply optional unidirectional updates mu^5 . A recover statement r is a unidirectional update of the form:

$$r ::= \text{if } e \text{ then } r \text{ else } r' \mid \text{delete} \mid \text{keep } u \\ \mid \text{case } e \text{ of } \vec{pat} \rightarrow \vec{r}$$

It supports conditionals and case expressions like regular updates, and two primitive operations: `delete`, to delete an unmatched source; and `keep u` , to keep an unmatched source modified with u so that it does not satisfy the above e_s filtering expression.

4.5 Core bidirectional lens language

To interpret our core language bidirectionally, we translate core bidirectional updates into a core language of “putback”-style lenses over generalized tree structures [27], supporting regular expression types. For the context of this paper, a *lens* is a BX $l : \tau_s \Leftrightarrow_{\Gamma} \tau_v$ between a source type τ_s and a view type τ_v under an environment of type Γ , defined using the combinators from Figure 10. The concrete syntax for lenses is not essential in our design, and can be skipped unless for understanding the bidirectional update semantics discussed in Section 5. The intuition for each combinator (read as a transformation from view to source) should be understandable from its type signature, and more details regarding their concrete bidirectional semantics can be found in [27]⁶.

Each lens in the above language comprises two partial functions $U : \Gamma \rightarrow \text{Maybe } \tau_s \rightarrow \tau_v \rightarrow \tau_s$ and $Q : \tau_s \rightarrow \tau_v$, satisfying laws similar to `UPDATEQUERY` and `QUERYUPDATE`. Since these functions are partial, updating or querying may fail at runtime. This is sometimes inevitable, for instance, whenever a view value does not satisfy a view condition or a view dependency written in a `WHERE VIEW` clause. Remark that the update function U receives an additional environment of type Γ and an optional source of type *Maybe* τ_s (as in Haskell, but with short-hand notation *Nothing* = . and *Just* $v = v$ for values), to account for cases (like `UNMATCHV` clauses) when a new source must be reconstructed from the view without updating an existing source [5, 26]. For a lens polymorphic over its environment type, we often write just $l : \tau_s \Leftrightarrow \tau_v$.

⁵ We often refer to optional updates mu , optional values mv , etc by defining grammars $mu ::= \cdot \mid u, mv ::= \cdot \mid v$ and so on.

⁶ An Haskell implementation of these (and more) combinators can be found at <http://hackage.haskell.org/package/putlenses>.

4.6 XML subtyping and ambiguity

The notion of subtyping plays a crucial role in XML approaches with regular expression types. A type τ_1 is said to be a *subtype* of τ_2 , written $\tau_1 <: \tau_2$, if the flat values belonging to τ_1 are also values of τ_2 , i.e., $\llbracket \tau_1 \rrbracket_{flat} \subseteq \llbracket \tau_2 \rrbracket_{flat}$. Under a flat value representation, a value of a type naturally belongs simultaneously to all its supertypes. This motivates a notion of structural equality between types: two types τ_1 and τ_2 are equivalent ($\tau_1 =: \tau_2$), meaning that they accept the same set of flat values, if both $\tau_1 <: \tau_2$ and $\tau_2 <: \tau_1$, e.g., $\tau =: \tau \mid \tau$. It also induces an equivalence relation \sim that ignores structure and relates values parsing the same data using different markup, formally, $v \sim v' \triangleq flat(v) = flat(v')$, e.g., $L v \sim R v$.

A type τ is said to be *unambiguous* if the equivalence relation for structured values of that type (\sim_{τ}) is the equality relation ($=_{\tau}$), intuitively meaning that there is only one way to parse a flat value of type τ into a structured value of type τ . Unambiguous regular expression types have a direct correspondence to algebraic data types [29], and standard automata algorithms exist for deciding unambiguity of regular expression types in polynomial time [10, 31].

Since we retain a structured representation of values, upcasting a value v_1 of type τ_1 into a supertype τ_2 requires more than a proof of subtyping: we must also change v_1 into a value v_2 that contains the same flat information as v_1 but conforms to the structure of τ_2 . This problem has been considered in [23], that introduces a subtyping algorithm as a proof system with judgments of the form $\vdash \tau_1 <: \tau_2 \Rightarrow c$, that we treat as a “black box”. In BX terms, $c : \tau_2 \Leftarrow \tau_1$ is called a *canonizer* [14], which is a bit like a lens from τ_2 to τ_1 that comprises a total upcast function $ucast : \tau_1 \rightarrow \tau_2$, and a partial downcast function $dcast : \tau_2 \rightarrow \tau_1$. In our sense, canonizers satisfy two properties stating that they only handle structure:

$$\begin{array}{l} ucast \ v_1 \sim v_1 \qquad \text{UP}\sim \\ dcast \ v_2 = v_1 \Rightarrow v_1 \sim v_2 \qquad \text{DOWN}\sim \end{array}$$

For an unambiguous type τ_2 , we can lift a canonizer $c : \tau_1 \Leftarrow \tau_2$ into a well-behaved lens lift $\hat{c} : \tau_1 \Leftrightarrow \tau_2$. Similarly, for an unambiguous type τ_1 , the inverse c^{-1} of a canonizer $c : \tau_1 \Leftarrow \tau_2$ —even though not a canonizer itself due to partiality—can be lifted into a well-behaved lens lift $\hat{c}^{-1} : \tau_2 \Leftrightarrow \tau_1$.

5. Type System and Staged Semantics

In traditional XML-based languages with regular expression types [8, 9, 17], the loose separation between types and values lends itself to a Curry-style interpretation: terms in the language are given semantics regardless of typing, and typing rules assess the well-formedness of terms. In BIFLUX, however, types come prior to semantics, as an update program describes a BX between two known source and view types on which its meaning may naturally depend. This leads to a Church-style interpretation, where semantics is given to type derivations instead of independent terms. Also, as we have seen before, updates in our core language are given semantics in two stages: they are first interpreted as lens expressions (typed between given types), and lenses have themselves a bidirectional semantics.

5.1 Interpreting expressions, paths and patterns

The type system and semantics for expressions and paths is similar to that of μXQ [9] and `FLUX` [8], adapted to our context. We merge typing and semantics into a single judgment $\Gamma; mx \vdash e : \tau \Rightarrow \lambda\gamma.v$, meaning “in type environment Γ with optional root variable mx , expression e has type τ , and given environment γ it evaluates to value v ”. Type soundness is guaranteed by construction, such that the argument environment γ is of type Γ and the produced v is of type τ . The concrete rules can be found in our technical report [28]. We also define a judgment $\Gamma; \tau \vdash e : \tau' \Rightarrow \lambda\gamma.v.v'$, meaning “in type environment Γ with root type τ , expression e has type τ' , and

$(\circ\lt)$	$:(\tau_1 \Leftrightarrow_{\Gamma} \tau_2) \rightarrow (\tau_2 \Leftrightarrow_{\Gamma} \tau_3) \rightarrow (\tau_1 \Leftrightarrow_{\Gamma} \tau_3)$	ifSthenelse	$:(\tau_1 \rightarrow \text{bool}) \rightarrow (\tau_1 \Leftrightarrow_{\Gamma} \tau_2) \rightarrow (\tau_1 \Leftrightarrow_{\Gamma} \tau_2) \rightarrow (\tau_1 \Leftrightarrow_{\Gamma} \tau_2)$
unfork	$:(\tau_1 \Leftrightarrow_{\Gamma} \tau_3) \rightarrow (\tau_2 \Leftrightarrow_{\Gamma} \tau_3) \rightarrow (\tau_1, \tau_2 \Leftrightarrow_{\Gamma} \tau_3)$	ifVthenelse	$:(\tau_2 \rightarrow \text{bool}) \rightarrow (\tau_1 \Leftrightarrow_{\Gamma} \tau_2) \rightarrow (\tau_1 \Leftrightarrow_{\Gamma} \tau_2) \rightarrow (\tau_1 \Leftrightarrow_{\Gamma} \tau_2)$
remfst	$:(\tau_2 \rightarrow \tau_1) \rightarrow (\tau_2 \Leftrightarrow_{\Gamma} (\tau_1, \tau_2))$	withEnv	$:(\text{Maybe } \tau_1 \rightarrow \tau_2 \rightarrow \Gamma \rightarrow \Gamma') \rightarrow (\tau_1 \Leftrightarrow_{\Gamma'} \tau_2) \rightarrow (\tau_1 \Leftrightarrow_{\Gamma} \tau_2)$
map	$:(\tau_1 \Leftrightarrow_{\Gamma} \tau_2) \rightarrow (\tau_1^* \Leftrightarrow_{\Gamma} \tau_2^*)$	keep	$:\tau \Leftrightarrow_{\Gamma} ()$
keepfst	$:(\tau_1, \tau_2) \Leftrightarrow_{\Gamma} \tau_2$	keepsnd	$:(\tau_1, \tau_2) \Leftrightarrow_{\Gamma} \tau_1$
alignpos	$:(\tau_1 \rightarrow \text{bool}) \rightarrow \text{Maybe } (\Gamma \rightarrow \tau_2 \rightarrow \tau_1) \rightarrow (\Gamma \rightarrow \tau_1 \rightarrow \text{Maybe } \tau_1) \rightarrow (\tau_1 \Leftrightarrow_{\Gamma} \tau_2) \rightarrow (\tau_1^* \Leftrightarrow_{\Gamma} \tau_2^*)$	bot	$:\tau_1 \Leftrightarrow_{\Gamma} \tau_2$
alignkey	$(\Gamma \rightarrow \tau_1 \rightarrow \tau_{k_1}) \rightarrow (\Gamma \rightarrow \tau_2 \rightarrow \tau_{k_2}) \rightarrow (\tau_{k_1} \rightarrow \tau_{k_2} \rightarrow \text{bool})$ $\rightarrow (\tau_1 \rightarrow \text{bool}) \rightarrow \text{Maybe } (\Gamma \rightarrow \tau_2 \rightarrow \tau_1) \rightarrow (\Gamma \rightarrow \tau_1 \rightarrow \text{Maybe } \tau_1) \rightarrow (\tau_1 \Leftrightarrow_{\Gamma} \tau_2) \rightarrow (\tau_1^* \Leftrightarrow_{\Gamma} \tau_2^*)$	id	$:\tau \Leftrightarrow_{\Gamma} \tau$
ifthenelse	$:(\Gamma \rightarrow \text{Maybe } \tau_1 \rightarrow \tau_2 \rightarrow \text{bool}) \rightarrow (\tau_1 \Leftrightarrow_{\Gamma} \tau_2) \rightarrow (\tau_1 \Leftrightarrow_{\Gamma} \tau_2) \rightarrow (\tau_1 \Leftrightarrow_{\Gamma} \tau_2)$	in	$:n[\tau] \Leftrightarrow_{\Gamma} \tau$
		listeq	$:\tau^* \Leftrightarrow_{\Gamma} \tau$

Figure 10: Language of point-free lenses for translating core bidirectional updates.

given environment $\gamma : \Gamma$ and root value $v : \tau$ it evaluates to value $v' : \tau'$. The root type and underlying value are added to the (type and value) environments under a fresh root variable:

$$\frac{\Gamma[x : \tau]; x \vdash e : \tau' \Rightarrow \lambda v \gamma[x := v].v' \quad x \notin \text{dom}(\Gamma)}{\Gamma; \tau \vdash e : \tau' \Rightarrow \lambda \gamma v.v'}$$

In BiFLUX, we follow a simple approach to pattern matching. We first define pattern type inference as a judgment $\vdash \text{pat} : \tau \Rightarrow \Pi_{\tau}$ that reads “pattern *pat* has type τ and yields a lens environment Π_{τ} ”. This is straightforward since we require all pattern variables to be annotated with a type; many others, including [16, 32], have studied more advanced XML-based pattern type inference techniques. A *lens environment* Π_{τ} is a set of bindings $x := (l, \tau')$ of variables to lenses from a source type τ , such that $l : \tau \Leftrightarrow \tau'$. As an abuse of notation, we see Π_{τ} as a type environment; we cast a source value v of type τ into an environment of source variables $\gamma_{v, \Pi_{\tau}} = \{x_i := Q_i(v) \mid x_i := (l_i, \tau_i) \in \Pi_{\tau}\}$.

The second step for matching a pattern against an input value is to simply check for subtyping between the input type and the inferred pattern type [16], as the witness functions of our subtyping algorithm already give us a way to convert values in both directions. For ambiguous patterns that may match an input value in multiple ways, such functions are responsible for resolving the ambiguity. Different matching policies can be modelled by adapting the subtyping proof system, as noted in [23], despite not being essential for our approach.

5.2 Interpreting unidirectional updates

The typechecking and operational semantics of core unidirectional updates mimic those of the core FLUX language, and can be found in the companion technical report [28]. The judgment $\Gamma \vdash \{\tau\} u \{\tau'\} \Rightarrow \lambda \gamma v.v'$ means that “in type environment Γ , a unidirectional update maps values of type τ to values of type τ' , and given environment γ and input value v (of type τ) it produces the updated value v' (of type τ')”.

5.3 Interpreting bidirectional updates

Contrarily to unidirectional updates, that modify values from an input to an output type, bidirectional updates are evaluated against two given source and view types and return a lens between those types. The judgment $\Gamma; \Pi_{\tau} \vdash \{\tau\} b \{\nu\} \Rightarrow l$ indicates that “in type environment Γ and lens environment Π_{τ} , a statement b produces a lens l (between source type τ and view type ν under type environment Γ)”. The type environment Γ denotes normal variables (that are in scope of the update function of the generated lens); the lens environment Π_{τ} denotes source variables from the source τ ; and the view type ν defines a view environment (also referred to as ν) that denotes view variables from the view ν . The dichotomy between our representations of source and view environments is justified by the need to keep a special account of view information: declared view variables must be used at least once in the update,

while source variables might not be used at all to update the source. Most rules for typechecking and evaluating bidirectional updates as lenses are shown in Figure 11. (The *listifyS* and *listifyV* functions and the judgment $\vdash \tau_1 \leq \tau_2 \Rightarrow l$ will be introduced in Section 5.4.)

Source/view expressions are interpreted under only a lens/view environment: “in lens environment Π_{τ} (and source type τ), a source expression e has type τ' , and given a value $v : \tau$ produces a value $v' : \tau'$ ”; “in view environment ν (and view type ν), a view expression e has type τ' , and given a value $v : \nu$ produces a value $v' : \tau'$ ”.

$$\frac{\Gamma; x \vdash e : \tau' \Rightarrow \lambda \gamma_{v, \Pi_{\tau}}[x := v].v' \quad \Gamma = \Pi_{\tau}[x : \tau] \quad x \notin \text{vars}(\Pi_{\tau})}{\Pi_{\tau} \vdash_S e : \tau' \Rightarrow \lambda v.v'} \quad \frac{\Gamma; \nu; \cdot \vdash e : \tau' \Rightarrow \lambda \gamma_{v: \nu}.v'}{\nu \vdash_V e : \tau' \Rightarrow \lambda v.v'}$$

Basic combinators The *skip* combinator returns the lens that does not update the source if the view is the empty sequence, while *fail* yields the bottom lens that always fails to update the source; *replace* completely replaces the source with a view that is “smaller”, such that the view type is a subtype of the source type.

Composition Bidirectional composition $b_1; b_2$ first splits the view into two parts ν_1 and ν_2 , evaluating b_1 with view ν_1 followed by b_2 with view ν_2 over the same source. The auxiliary function $\text{split}(\nu, \theta_1, \theta_2) = (\nu_1, \nu_2, l_{12})$ splits ν such that $\text{dom}(\nu_1) = \text{dom}(\nu) \cap \theta_1$ and $\text{dom}(\nu_2) = \text{dom}(\nu) \cap \theta_2$, returning a lens $l_{12} : (\nu_1, \nu_2) \Leftrightarrow \nu$. Parallel lens composition ($\text{unfork } l_1 \ l_2$) is different from sequential lens composition ($l_1 \circ l_2$); for it to be well-behaved, the second update can not affect the result of the first query, i.e., $Q_1(U_2(v, v_2)) \sqsubseteq Q_1(v)$. This property is currently checked dynamically by *unfork*, that fails to update the source otherwise. We believe that this check could be done at static time by combining recent work on XML query-update independence [2, 4, 6].

Changing source focus The bidirectional update $p[b]$ changes the source focus by traversing down the source path p , and then evaluates b with a fresh lens environment. Note that it is intrinsically different from the unidirectional update $p[u]$, as the source type does not change and the source data is not updated in-place, but modified to embed the view data. Precisely, this requires interpreting the source path as a lens via a judgment $\Pi_{\tau} \vdash_S \{\tau\} p \{\tau'\} \Rightarrow l$, that reads “in lens environment Π_{τ} , path p changes the source focus from type τ to type τ' , and produces a lens $l : \tau \Leftrightarrow \tau'$ ”. Like in FLUX, arbitrary paths can not be used to change the source focus (see [28]). For example, only the *self* and *child* axes (and no absolute paths) are supported; this ensures that only descendants of the source focus can be selected as the new source focus and that a selection contains no overlapping elements. The *iter* b operator shifts the source focus to all tree values in a source forest, and runs b for each tree. Its corresponding BX will duplicate the view value for each source tree during updates, and enforce all selected source trees to be the same during queries.

$$\boxed{\Gamma; \Pi \vdash \{\tau\} s \{\nu\} \Rightarrow l}$$

$$\begin{array}{c}
\frac{}{\Gamma; \Pi_\tau \vdash \{\tau\} \text{skip } \{()\} \Rightarrow \text{keep}} \quad \frac{}{\Gamma; \Pi_\tau \vdash \{\tau\} \text{fail } \{\nu\} \Rightarrow \text{bot}} \\
\frac{\Gamma; \Pi_\tau \vdash \{\tau\} b_1 \{\nu_1\} \Rightarrow l_1 \quad \Gamma; \Pi_\tau \vdash \{\tau\} b_2 \{\nu_2\} \Rightarrow l_2 \quad \text{split}(\nu, \text{vars}(b_1), \text{vars}(b_2)) = (\nu_1, \nu_2, l_{12})}{\Gamma; \Pi_\tau \vdash \{\tau\} b_1; b_2 \{\nu\} \Rightarrow \text{unfork } l_1 \ l_2 \circ l_{12}} \quad \frac{\Gamma; \Pi_\tau \vdash \{\tau\} b_1 \{\nu\} \Rightarrow l_1 \quad \Gamma; \Pi_\tau \vdash \{\tau\} b_2 \{\nu\} \Rightarrow l_2 \quad \Pi_\tau \vdash_S e : \text{bool} \Rightarrow f}{\Gamma; \Pi_\tau \vdash \{\tau\} \text{ifS } e \text{ then } b_1 \text{ else } b_2 \{\nu\} \Rightarrow \text{ifSthenelse } f \ l_1 \ l_2} \\
\frac{\Pi_\tau \vdash_S \{\tau\} p \ \{\tau'\} \Rightarrow l_1 \quad \Gamma; \emptyset \vdash \{\tau'\} b \ \{\nu\} \Rightarrow l_2}{\Gamma; \Pi_\tau \vdash \{\tau\} p[b] \ \{\nu\} \Rightarrow l_1 \circ l_2} \quad \frac{\text{listifyS}(\tau) = (\tau_1, l_1) \quad \Gamma; \emptyset \vdash \{\tau_1\} b \ \{\nu\} \Rightarrow l_2}{\Gamma; \Pi_\tau \vdash \{\tau\} \text{iter } b \ \{\nu\} \Rightarrow l_1 \circ \text{map } l_2 \circ \text{listeq}} \\
\frac{\Gamma[x_V : \tau']; \Pi_\tau \vdash \{\tau\} b \ \{x_V : \tau'\} \Rightarrow l_1 \quad \Gamma \vdash_V \{\tau'\} e \ \{\nu\} \Rightarrow l_2 \quad f_1 \ m_V \ s \ v'_o \ \gamma' = \gamma[x_V := v'_o]}{\Gamma; \Pi_\tau \vdash \{\tau\} [b]e \ \{\nu\} \Rightarrow \text{withEnv } f_1 \ l_1 \circ l_2} \quad \frac{\Gamma; \Pi_\tau \vdash \{\tau\} b \ \{\nu_2\} \Rightarrow l_2 \quad \nu_2 \vdash_V e : \nu_1 \Rightarrow f_{21} \quad x \in \text{dom}(\nu) \quad \text{split}(\nu, \{x\}, \text{dom}(\nu) \setminus \{x\}) = (\nu_1, \nu_2, l_{12})}{\Gamma; \Pi_\tau \vdash \{\tau\} \text{view } x := e \ \text{in } b \ \{\nu\} \Rightarrow l_2 \circ \text{remfst } f_{21} \circ l_{12}} \\
\frac{\text{listifyS}(\tau) = (\tau_1, l_1) \quad \text{listifyV}(\tau') = (\tau_2, l_2) \quad \Gamma[x_V : \tau_2]; \emptyset \vdash \{\tau_1\} b \ \{x_V : \tau_2\} \Rightarrow l \quad \emptyset \vdash_S e : \text{bool} \Rightarrow f_e \quad \Gamma \vdash_{\text{create}} \{\tau_2\} c \ \{\tau_1\} \Rightarrow m_f \ c \quad \Gamma \vdash_{\text{recover}} r \ \{\tau\} \Rightarrow f_r \quad f \ m_V \ s \ v_v \ \gamma = \gamma[x_V := v_v]}{\Gamma; \Pi_\tau \vdash \{\tau\} \text{alignpos } e \ b \ c \ r \ \{x_V : \tau'\} \Rightarrow l_1 \circ \text{alignpos } f_e \ m_f \ c \ f_r \ (\text{withEnv } f \ l) \circ l_2} \\
\frac{\Gamma; \Pi_\tau \vdash \{\tau\} b_1 \{\nu\} \Rightarrow l_1 \quad \Gamma; \Pi_\tau \vdash \{\tau\} b_2 \{\nu\} \Rightarrow l_2 \quad \nu \vdash_V e : \text{bool} \Rightarrow f}{\Gamma; \Pi_\tau \vdash \{\tau\} \text{ifV } e \text{ then } b_1 \text{ else } b_2 \{\nu\} \Rightarrow \text{ifVthenelse } f \ l_1 \ l_2} \quad \frac{}{\Gamma; \Pi_\tau \vdash \{\tau\} \text{replace } \{\nu\} \Rightarrow l} \\
\frac{\Gamma; \Pi_\tau \vdash \{\tau\} b_1 \{\nu\} \Rightarrow l_1 \quad \Gamma; \Pi_\tau \vdash \{\tau\} b_2 \{\nu\} \Rightarrow l_2 \quad \Gamma; \tau \vdash e : \text{bool} \Rightarrow \lambda \gamma \ v_s. b \quad f \cdot v_v = \text{true} \quad f \ v_s \ v_v \ \gamma = b}{\Gamma; \Pi_\tau \vdash \{\tau\} \text{if } e \text{ then } b_1 \text{ else } b_2 \{\nu\} \Rightarrow \text{ifthenelse } f \ l_1 \ l_2} \\
\frac{P(\vec{x} : \vec{\tau}) : \nu_s \Leftrightarrow \nu_v \in \Delta \quad \Pi_\tau \vdash_{\text{procS}} \{\tau\} \vec{p}_s \ \{\nu_s\} \Rightarrow l_s \quad \Gamma \vdash_{\text{procV}} \{\nu_v\} \vec{e}_v \ \{\nu\} \Rightarrow l_v \quad \Gamma; \tau \vdash e : \tau_1 \Rightarrow \lambda \gamma \ v. v_1 \quad \dots \quad \Gamma; \tau \vdash e : \tau_n \Rightarrow \lambda \gamma \ v. v_n \quad f \ v_s \ v_v \ \gamma = \{x_1 := v_1, \dots, x_n := v_n\}}{\Gamma; \Pi_\tau \vdash \{\tau\} P(\vec{p}_s, \vec{e}_v, \vec{e}) \ \{\nu\} \Rightarrow l_s \circ \text{withEnv } f \ P \circ l_v}
\end{array}$$

$$\boxed{\Gamma \vdash_{\text{create}} \{\tau_2\} c \ \{\tau_1\} \Rightarrow \lambda_m \gamma \ v_2. v_1}$$

$$\boxed{\Gamma \vdash_{\text{recover}} r \ \{\tau\} \Rightarrow \lambda \gamma \ v. mv}$$

$$\begin{array}{c}
\frac{\Gamma[x_V : \tau] \vdash \{\tau\} u \ \{\tau'\} \Rightarrow \lambda \gamma [x_V := v] \ v. v'}{\Gamma \vdash_{\text{create}} \{\tau\} u \ \{\tau'\} \Rightarrow \lambda \gamma \ v. v'} \quad \frac{\Gamma \vdash \{\tau\} u \ \{\tau\} \Rightarrow \lambda \gamma \ v. v'}{\Gamma \vdash_{\text{recover}} \text{keep } u \ \{\tau\} \Rightarrow \lambda \gamma \ v. v'} \quad \frac{}{\Gamma \vdash_{\text{recover}} \text{delete } u \ \{\tau\} \Rightarrow \lambda \gamma \ v. \cdot} \\
\frac{\Gamma[x : \tau]; x \vdash e : \text{bool} \Rightarrow \lambda \gamma [x := v]. \text{true} \quad \Gamma \vdash_{\text{recover}} r \ \{\tau\} \Rightarrow \lambda \gamma \ v. mv}{\Gamma \vdash_{\text{create}} \{\tau\} \cdot \ \{\tau'\} \Rightarrow \cdot} \quad \frac{\Gamma[x : \tau]; x \vdash e : \text{bool} \Rightarrow \lambda \gamma [x := v]. \text{false} \quad \Gamma \vdash_{\text{recover}} r' \ \{\tau\} \Rightarrow \lambda \gamma \ v. mv}{\Gamma \vdash_{\text{recover}} \text{if } e \text{ then } r \ \text{else } r' \ \{\tau\} \Rightarrow \lambda \gamma \ v. mv} \\
\frac{}{\Gamma \vdash_{\text{create}} \{\tau\} \cdot \ \{\tau'\} \Rightarrow \cdot} \quad \frac{}{\Gamma \vdash_{\text{recover}} \text{if } e \text{ then } r \ \text{else } r' \ \{\tau\} \Rightarrow \lambda \gamma \ v. mv}
\end{array}$$

Figure 11: Bidirectional update well-formedness and semantics.

Changing view focus The operation $[b]e$ changes the view focus according to the non-source expression e , and then evaluates b for the intermediate view⁷. The judgment $\Gamma \vdash_V \{\tau\} e \ \{\nu\} \Rightarrow l$ indicates that “in type environment Γ , expression e changes the view focus from type ν to type τ , and produces a lens $l : \tau \Leftrightarrow_{\Gamma} \nu$ ” [28]. Note that, to be successfully interpreted as a “backward” lens, the expression e may only add information to the view, since all view information must be used to update the source; this implies that only a restricted subset of our path language that can be statically checked to be injective [28] can be used to change the view focus. The operation $\text{view } x := e \ \text{in } b$ removes a view variable x from the current view environment without loss of information (with the view expression e evidencing how x can be computed from the other view variables), followed by running b .

Conditionals The conditional operations ifS , ifV and if choose between two statements b_1 or b_2 according to a boolean expression e , and differ subtly on the bidirectional behavior of the underlying lenses (see [27]), given that e is a source, view or normal expression. The same rationale is applied to case expressions [28]).

Source-view alignment The alignment operations alignpos and alignkey synchronize source and view forests. They 1) uniformize the source and view types into lists using listifyS and listifyV , 2) align source and view elements using matching algorithms [1] on lists by position or by keys —calculated as paths on the current source/view— and 3) run the statement b for matching source/view pairs, the create statement c for unmatched source elements or the recover statement r for unmatched source elements. Create statements are interpreted unidirectionally using a judgment $\Gamma \vdash_{\text{create}} \{\tau_2\} c \ \{\tau_1\} \Rightarrow \lambda_m \gamma \ v_2. v_1$, saying that “in type environment Γ , a create statement c between source type τ_1 and view type τ_2 returns an optional function (denoted by $\lambda_m \cdot \cdot$) that given an environment $\gamma : \Gamma$ and a view value v_2 creates a source value v_1 ”. In reverse direction, the judgment $\Gamma \vdash_{\text{recover}} r \ \{\tau_1\} \Rightarrow \lambda \gamma \ v_1. mv$ states that “in type environment Γ , a recover statement r for a source type τ_1 returns a function that given an environment $\gamma : \Gamma$ and a source value $v_1 : \tau_1$ returns an optional recovered source value $mv_1 : \text{Maybe } \tau_1$ ”. Note that our alignment operations receive a source predicate e denoting which source values have a correspondence to view values; the underlying lenses enforce that newly created source values must satisfy e , while recovered source values do not originate from the view and must not satisfy e .

⁷We use x_V as a special internal view variable to accommodate the intermediate view as a record type.

Procedures Procedure calls $P(\vec{p}_s, \vec{e}_v, \vec{e})$ are interpreted under a new type environment computed from the argument expressions \vec{e} . The new lens environment is computed from the source paths \vec{p}_s , via a judgment $\Pi_\tau \vdash_{\text{procs}} \{\tau\} \vec{p} \{\nu_s\} \Rightarrow l$, and the new view environment is computed from the non-source expressions \vec{e}_v via a judgment $\Gamma \vdash_{\text{procv}} \{\nu_v\} \vec{e} \{\nu\} \Rightarrow l$. The semantics of a BiFluX program is given by converting a set Δ of procedure declarations $P(\vec{x} : \vec{\tau}) : \nu_s \Leftrightarrow \nu_v \triangleq s$ into a set Δ_{\Leftrightarrow} of lenses $P = l : \nu_s \Leftrightarrow \{\tau_1, \dots, \tau_n : \tau_n\} \nu_v$, according to the following derivation:

$$\frac{\begin{array}{c} \{x_1 : \tau_1, \dots, x_n : \tau_n\} \cup \nu_s \cup \nu_v; \Pi_{\nu_s} \vdash \{\nu_s\} s \{\nu_v\} \Rightarrow l \\ f \nu_s \nu_v \gamma = \gamma \cup \gamma_{\nu_s : \nu_s} \cup \gamma_{\nu_v : \nu_v} \end{array}}{\vdash P(\vec{x} : \vec{\tau}) : \nu_s \Leftrightarrow \nu_v \triangleq s \Rightarrow P = \text{withEnv } f \ l}$$

5.4 Type normalization

The semantics of our core language relies instrumentally on subtyping to match source and view types. The simplest example is the `replace` bidirectional update, that requires the view type τ_2 to be “smaller” than the source type τ_1 (Figure 11), according to a judgment $\vdash \tau_2 \leq \tau_1 \Rightarrow l$ that returns as evidence a lens $l : \tau_1 \Leftrightarrow \tau_2$. As the reader may guess, we could compute $\vdash \tau_2 \leq \tau_1 \Rightarrow l$ by checking for subtyping between τ_2 and τ_1 ($\vdash \tau_2 <: \tau_1 \Rightarrow c$) and lifting the resulting canonizer into a lens. But, before doing so, we must verify that: the view type τ_2 is unambiguous, a requirement to lift c into a lens $\text{lift } c$; and the source type τ_1 is unambiguous, since the `ucast` function of the canonizer does not consider the original source and a naive implementation of an update would potentially discard source information projected away by the lens.

As an example, recapitulate the source database of books used in Section 2 and imagine how we could evaluate a simple update:

REPLACE \$source/books/book/title **WITH** \$view

using a list `title[string]*` as the view type. Consistently with the unidirectional semantics for paths [28], the evaluation of the source path as a lens traverses down the source tree and keeps only titles, destroying element labels and replacing all authors for the empty sequence, and modifies the source focus to the intermediate type `(title[string], (), ())*`. Then `replace` checks for subtyping between the view type and the intermediate type, producing a mediating canonizer. Since the intermediate type is ambiguous, the upcasting function of the canonizer would translate a view sequence `[title["mybook"]]` into an intermediate value `[(title["mybook"], ((), []))]`, that would in turn lead to an updated source `books[[book[(title["mybook"], (aut1, []))]]]`. This is clearly unsatisfactory as it would discard the entire book database except the first author `aut1` of the first book!

Source and view normalization To avoid these problems, we introduce *type normalization* procedures that simplify regular expression types into unambiguous ones while carefully preserving the original markup information that keeps trace of hidden information. Namely, we define a *source normalization* procedure $\text{normS}(\tau) = (\tau', l)$, that normalizes a source type into an unambiguous subtype τ' together with a lens $l : \tau \Leftrightarrow \tau'$, and a *view normalization* procedure $\text{normV}(\tau) = (\tau', l)$, that normalizes a view type into an unambiguous subtype type τ' together with a lens $l : \tau \Leftrightarrow \tau'$. Their main difference is that they compute lenses in opposite directions, suggesting that source normalization may abstract ambiguous information (like redundant choices) while view normalization may not [28]. We do not claim that our normalization procedures are complete, in the sense that they can disambiguate any ambiguous type, but when they (statically) do succeed the normalized types are unambiguous. In a nutshell, we try to normalize a source type using automata reduction techniques [10, 31], and derive a lens between the ambiguous and unambiguous types; we

only normalize view types into isomorphic types. Furthermore, we evaluate source paths in a two-phased semantics: a path is interpreted as a completely information-preserving lens, that just marks unselected source types τ with a special tag $\boxed{\tau}$ instead of replacing them for the empty sequence; and tagged types are to be removed during source normalization, that already requires special handling anyway.

Normalized subtyping and listification We are now ready to safely define $\tau_1 \leq \tau_2$ as the following judgment:

$$\frac{\text{normS}(\tau_2) = (\tau'_2, l_2) \quad \text{normV}(\tau_1) = (\tau'_1, l_1) \quad \vdash \tau'_1 <: \tau'_2 \Rightarrow c}{\vdash \tau_1 \leq \tau_2 \Rightarrow l_1 \circ \text{lift } c \circ l_2}$$

Unlike iteration for unidirectional updates, that changes the focus to all the atomic elements in a forest and updates their values (and types!) independently, bidirectional iteration is type-preserving and involves updating a source sequence with view information and somehow fitting it back into a shape that conforms to the source type. Since *shape alignment* is an intractable problem for arbitrary source and view types [26], we uniformize source and view types into lists of (choices of) atomic types according to two functions:

$$\frac{\begin{array}{c} \text{elems}(\tau') = \{\alpha_1, \dots, \alpha_n\} \quad (\alpha_0 \mid \dots \mid \alpha_n)^* \text{ unambiguous} \\ \text{normS}(\tau) = (\tau', l) \quad \vdash \tau' <: (\alpha_0 \mid \dots \mid \alpha_n)^* \Rightarrow c \end{array}}{\text{listifyS}(\tau) = ((\alpha_0 \mid \dots \mid \alpha_n), l \circ \text{lift } c^{-1})}$$

$$\frac{\begin{array}{c} \text{elems}(\tau') = \{\alpha_0, \dots, \alpha_n\} \quad (\alpha_0 \mid \dots \mid \alpha_n)^* \text{ unambiguous} \\ \text{normV}(\tau) = (\tau', l) \quad \vdash \tau' <: (\alpha_0 \mid \dots \mid \alpha_n)^* \Rightarrow c \end{array}}{\text{listifyV}(\tau) = ((\alpha_0 \mid \dots \mid \alpha_n), \text{lift } c \circ l)}$$

Note that the uniformized list types may be more flexible supertypes, e.g. $(\alpha_1, \alpha_2)^* <: (\alpha_1 \mid \alpha_2)^*$, allowing more values than those that fit the real type. While this does not pose a problem with the for-each semantics of plural high-level updates, it may lead to runtime errors for **UPDATE FOR VIEW** updates over intricate structures. An alternative is to statically check for type equivalence in `listifyS`, supporting only pure source lists. The function $\text{elems} : \text{Type} \rightarrow \{\text{Atom}\}$ returns the set of atomic types in a sequence type.

6. BIFLUX to Core Update Normalization

In this section, we formalize the translation from the high-level BiFluX language to the core language presented in Section 4, highlighting the significant gap between them. This process is usually referred to as *normalization* in languages like XQuery and FLUX. We define two main normalization functions that interpret statements as bidirectional $\llbracket - \rrbracket_{\text{Stmt}}^b$ and unidirectional updates $\llbracket - \rrbracket_{\text{Stmt}}^u$. Most translation rules are straightforward and a few interesting ones are shown in Figure 12; the complete set can be found in [28]. Simple bidirectional updates are translated by a function $\llbracket - \rrbracket_{\text{Upd}}^b(e_s, e_v, \vec{x} = \vec{e})$, where the extra parameters group the **WHERE** clauses of the update into a source selection expression e_s , a view selection expression e_v and a sequence of view bindings $\vec{x} = \vec{e}$; these triples are parsed from a set of conditions, according to their source/view tags, by a function $\llbracket - \rrbracket_{\text{Conds}}^c$. Simple unidirectional updates are translated by a function $\llbracket - \rrbracket_{\text{Upd}}^u(e)$, where e is the conjunction of all the **WHERE** clauses of the update.

For special **UPDATE FOR VIEW** statements, the `splitVStmt` function parses a `VStmt` into a matching statement and two optional unmatched-view and unmatched-source statements. Optional unmatched-view statements are translated using a function $\llbracket - \rrbracket_{\text{MStmt}}^c(\text{mpat})$ that takes an extra optional view pattern and returns a core create update; if no **UNMATCHV** clause is defined, the `U` function of the underlying lens will be evaluated without an original source. Optional unmatched-source statements are translated using a function $\llbracket - \rrbracket_{\text{MStmt}}^r(\text{mpat})$ that takes an extra optional source pattern and returns a core recover update; if no **UNMATCHS** clause is defined, all unmatched source elements are deleted by default.

The translation denotes a partial function from high-level BiFLUX to core BiFLUX. For example, `INSERT` is not supported for bidirectional updates, `UPDATE FOR VIEW` is not supported for unidirectional updates, and `CREATE` or `KEEP` are only supported under `UNMATCHV` or `UNMATCHS` clauses, respectively. We assume that paths and expressions are expressed in terms of our core languages; this is standard practice as normalization of XQuery expressions or XPath paths can be done independently. To simplify the presentation, we also assume explicit `SOURCE` and `VIEW` tags, though our implementation is elaborated to implicitly distinguish between source/view/normal expressions, using the additional environment information available at the time of typechecking the core language.

7. Related Work

XML update languages Several XML update languages have been proposed, including (among many others) XQuery! [15], FLUX [8] and the standard W3C XQuery Update Facility [30]. Even though the specification style, expressiveness and semantics of the XML updates that can be written may vary significantly, they all focus on updating XML documents in-place, i.e., updating selected parts of an XML document, keeping the remaining parts of the document unchanged. This means that update programs can be seen as unidirectional transformations that insert, delete or replace elements in a source document and produce an updated document conforming to a new target type. XML Updates in BiFLUX are different in that they determine how to update a source document (using some view information) while preserving its source type. This poses different (BX-related) challenges on how to deal with non-in-place updates (like `UPDATE FOR VIEW` statements that may change the cardinality of a sequence instead of updating each element), and therefore how to modify the remaining parts of the source document (e.g., by changing branching decisions) to accommodate the new data so that the updated source fits into the same type.

XML view updating In [12], the author studies the problem of updating XML views of relational databases by translating view updates written in the XQuery Update Facility into embedded SQL updates. The work of [22] supports updatable views of XML data by giving a bidirectional semantics to the XQuery Core language. The semantic bidirectionalization technique of [24] interprets various XQuery use cases as BXs by encoding them as polymorphic Haskell functions. The Multifocal language [25] allows writing high-level generic XML views that can be applied to multiple XML schemas, producing a view schema and a lens conforming to the schemas. In the four approaches, the programmer writes a view function and the system derives a suitable view update translation strategy using built-in techniques that he can not configure. In BiFLUX, he writes an update translation strategy directly as an update (over the source) and the system derives the uniquely related query.

XML bidirectional languages Many bidirectional programming languages support tree-structured or XML data formats. Two popular XML bidirectional languages are XSugar [7] and biXid [21], that describe XML-to-ASCII and XML-to-XML mappings as pairs of intertwined grammars. While XSugar restricts itself to bijective grammars, biXid considers nondeterministic specifications and BXs are inherently ambiguous. Most functional bidirectional programming languages are based on lenses [13, 18, 26, 27], and follow a combinatorial style that puts special emphasis on building complex lenses by composition of smaller combinators. Depending on the choice of combinators, lens languages can become very powerful at specifying application-specific behavior [1, 26, 27]. However, their lower-level nature also induces a more cumbersome programming style that makes it impractical and often unintuitive for users to build non-trivial BXs by piping together several small, surgical steps.

BiFLUX features a new programming by update paradigm, that enables the high-level syntax of relational languages such as XSugar and biXid while providing a handful of intuitive update strategies. Remember the huge gap between our high-level BiFLUX language (variables, procedures) and the core lens language that gives it semantics (canonical “point-free” combinators). In [19], we have proposed a simple treeless functional language for writing total *put* (or update) functions, such that existence of a well-behaved lens can be checked statically, and corresponding total *get* (or query) functions can be derived automatically. The most significant and innovative difference in BiFLUX is again the declarative surface language used to specify BXs as bidirectional update programs, at a notably higher-level of abstraction than native *put* functions.

Quotient lenses [14] propose loosening lenses modulo equivalence relations, for easing the processing of ad-hoc data formats. To enable compositional reasoning, most interesting quotient lens combinators (concatenation, union, iteration) require the equivalence relations to be decomposable. At first glance, we could lift our core lens language into quotient lenses to allow seamless composition with (subtyping) canonizers. However, since our notion of equivalence is not decomposable for ambiguous types, this would not overcome the need for our type normalization procedures.

8. Conclusions

In this paper, we propose a novel *bidirectional programming by update* paradigm that comes to light from the idea of extending a traditional update language with bidirectional features. Under the new paradigm, programmers write bidirectional updates that specify how to update a source document to reflect additional view information. We substantiate with examples that this enjoys a better tradeoff between the expressiveness and declarativeness of the written bidirectional programs, by allowing users to write directly, in a friendly notation and at a nice level of abstraction, a view update translation strategy that bundles all the pieces to build a BX.

To demonstrate the potential of this paradigm, we designed BiFLUX, a type-safe high-level bidirectional XML update language. We have fully implemented BiFLUX in Haskell and the code, together with additional examples, is available from the project’s website⁸. Our tool translates source and view DTDs into Haskell type declarations using the `HaXml` package⁹, and interprets bidirectional updates as bidirectional lens transformations between the given schemas in a robust manner: the implementation of the core-to-lens translation is strongly-typed, serving as a proof of soundness that helped us catching early programming errors at compile-time; and the bidirectional semantics is completely guided by an underlying lens language, whose correctness has been shown separately in [27]. To support a flexible design, with arbitrary conditionals, case statements and expressions, the statically generated lenses may undergo runtime checks (for particular source and view documents) to ensure that the underlying update and query functions are well-behaved.

As future work, we plan to provide more static guarantees to BiFLUX by incorporating existing path-query static analyses, implement more powerful pattern type inference algorithms to avoid excessive annotations, and extend the class of bidirectional updates that can be written by integrating user-defined lenses for defining source and view foci. We also plan to improve the efficiency of our prototype for large XML databases by exploring optimizations to the underlying lens language, including incremental update translation. To empirically study the practical impact of BiFLUX, we are currently undergoing a larger model-based code testing use case.

⁸ <http://www.prg.nii.ac.jp/projects/BiFlux>

⁹ <http://hackage.haskell.org/package/HaXml>

$$\begin{aligned}
\llbracket \text{pat IN } p \rrbracket_{PatPath}^{iter}(es, b) &= (p / \text{where } (\text{let } pat = \text{self in } es))[\text{iter } (\text{caseS self of } pat \rightarrow b)] \\
\llbracket \text{pat IN } p \rrbracket_{PatPath}^S(es, b) &= (p / \text{where } (\text{let } pat = \text{self in } es))[\text{caseS self of } pat \rightarrow b] \\
\llbracket \text{DELETE } patp \rrbracket_{Upd}^b(es, \text{true}, \cdot) &= \llbracket patp \rrbracket_{PatPath}^S(es, [\text{replace}]()) & \llbracket u \text{ WHERE } cs \rrbracket_{Stmt}^b = \llbracket u \rrbracket_{Upd}^b(\llbracket cs \rrbracket_{Conds}) \\
\llbracket \text{DELETE FROM } patp \rrbracket_{Upd}^b(es, \text{true}, \cdot) &= \llbracket patp \rrbracket_{PatPath}^S(es, \text{child}[\llbracket \text{replace} \rrbracket()]) & \llbracket s; s' \rrbracket_{Stmt}^b = \llbracket s \rrbracket_{Stmt}^b; \llbracket s' \rrbracket_{Stmt}^b \\
\llbracket \text{REPLACE } patp \text{ IN } p \text{ WITH } e' \rrbracket_{Upd}^b(es, e_V, \vec{x} = \vec{e}) &= \llbracket patp \rrbracket_{PatPath}^{iter}(es, \text{view } \vec{x} := \vec{e} \text{ in ifV } e_V \text{ then } [\text{replace}]e' \text{ else fail}) \\
\llbracket \text{REPLACE IN } patp \text{ WITH } e' \rrbracket_{Upd}^b(es, e_V, \vec{x} = \vec{e}) &= \llbracket patp \rrbracket_{PatPath}^{iter}(es, \text{view } \vec{x} := \vec{e} \text{ in ifV } e_V \text{ then child}[\llbracket \text{replace} \rrbracket]e' \text{ else fail}) \\
\llbracket \text{UPDATE } patp \text{ BY } s \rrbracket_{Upd}^b(es, e_V, \vec{x} = \vec{e}) &= \llbracket patp \rrbracket_{PatPath}^{iter}(es, \text{view } \vec{x} := \vec{e} \text{ in ifV } e_V \text{ then } \llbracket s \rrbracket_{Stmt}^b \text{ else fail}) \\
\llbracket \text{UPDATE } pat \text{ IN } p \text{ BY } vs \text{ FOR VIEW } pat' \text{ IN } p' \text{ MATCHING SOURCE BY } p_s \text{ VIEW BY } p_v \rrbracket_{Upd}^b(es, e_V, \vec{x} = \vec{e}) &= p[\llbracket b \rrbracket p'] \text{ where} \\
&((s_{SV}, ms_V, ms_S), p_s', p_v') = (\text{split } VStmt(vs), \text{case self of } pat \rightarrow p_s, \text{case self of } pat' \rightarrow p_v) \\
&b = \text{alignkey } (\text{case self of } pat \rightarrow es) p_s' p_v' \llbracket s_{SV} \rrbracket_{Stmt}^b \llbracket ms_V \rrbracket_{MStmt}^c(pat') \llbracket ms_S \rrbracket_{MStmt}^r(pat)
\end{aligned}$$

Figure 12: BiFluX bidirectional statement normalization.

Acknowledgments

This research was partially supported by the NSF under grant SHF-1253165 and by JSPS Grant-in-Aid for Scientific Research (A) No. 25240009.

References

- [1] D. M. J. Barbosa, J. Cretin, J. N. Foster, M. Greenberg, and B. C. Pierce. Matching lenses: alignment and view update. In *ICFP 2010*, pages 193–204. ACM, 2010.
- [2] M. Benedikt and J. Cheney. Destabilizers and independence of XML updates. *Proc. VLDB Endow.*, 3(1-2):906–917, 2010.
- [3] V. Benzaken, G. Castagna, and A. Frisch. CDuce: an XML-centric general-purpose language. In *ICFP 2003*, pages 51–63. ACM, 2003.
- [4] N. Bidoit, D. Colazzo, and F. Ulliana. Type-based detection of XML query-update independence. *Proc. VLDB Endow.*, 5(9):872–883, 2012.
- [5] A. Bohannon, J. N. Foster, B. C. Pierce, A. Pilkiewicz, and A. Schmitt. Boomerang: resourceful lenses for string data. In *POPL 2008*, pages 407–419. ACM, 2008.
- [6] S. Botcher. Testing intersection of XPath expressions under DTDs. In *IDEAS 2004*, pages 401–406. IEEE, 2004.
- [7] C. Brabrand, A. Möller, and M. I. Schwartzbach. Dual syntax for XML languages. *Information Systems*, 33(4-5):385–406, 2008.
- [8] J. Cheney. FLUX: functional updates for XML. In *ICFP 2008*, pages 3–14. ACM, 2008.
- [9] D. Colazzo, G. Ghelli, P. Manghi, and C. Sartiani. Static analysis for path correctness of XML queries. *Journal of Functional Programming*, 16(4-5):621–661, 2006.
- [10] H. Comon, M. Dauchet, R. Gilleron, C. Löding, F. Jacquemard, D. Lugiez, S. Tison, and M. Tommasi. *Tree automata techniques and applications*. Available at <http://www.grappa.univ-lille3.fr/tata>, 2007.
- [11] K. Czarniecki, J. N. Foster, Z. Hu, R. Lämmel, A. Schürr, and J. Terwilliger. Bidirectional transformations: A cross-discipline perspective. In *ICMT 2009*, volume 5563 of *LNCS*, pages 260–283. Springer, 2009.
- [12] L. Fegaras. Propagating updates through XML views using lineage tracing. In *ICDE 2010*, pages 309–320. IEEE, 2010.
- [13] J. N. Foster, M. B. Greenwald, J. T. Moore, B. C. Pierce, and A. Schmitt. Combinators for bidirectional tree transformations: A linguistic approach to the view-update problem. *ACM TOPLAS*, 29(3):17, 2007.
- [14] J. N. Foster, A. Pilkiewicz, and B. C. Pierce. Quotient lenses. In *ICFP 2008*, pages 383–396. ACM, 2008.
- [15] G. Ghelli, C. Ré, and J. Siméon. XQuery!: An XML query language with side effects. In *EDBT 2006*, volume 4254 of *LNCS*, pages 178–191. Springer, 2006.
- [16] H. Hosoya and B. Pierce. Regular Expression Pattern Matching for XML. In *POPL 2001*, pages 67–80. ACM, 2001.
- [17] H. Hosoya, J. Vouillon, and B. C. Pierce. Regular expression types for XML. In *ICFP 2000*, pages 11–22. ACM, 2000.
- [18] Z. Hu, S.-C. Mu, and M. Takeichi. A programmable editor for developing structured documents based on bidirectional transformations. *Higher-Order and Symbolic Computation*, 21(1-2):89–118, 2008.
- [19] Z. Hu, H. Pacheco, and S. Fischer. Validity checking of putback transformations in bidirectional programming (invited paper). In *FM 2014*, pages 1–15. Springer, 2014.
- [20] J. Hughes. Generalising monads to arrows. *Science of computer programming*, 37(1):67–111, 2000.
- [21] S. Kawanaka and H. Hosoya. biXid: a bidirectional transformation language for XML. In *ICFP 2006*, pages 201–214. ACM, 2006.
- [22] D. Liu, Z. Hu, and M. Takeichi. An expressive bidirectional transformation language for XQuery view update. *Progress in Informatics*, 10: 89–130, 2013.
- [23] K. Z. M. Lu and M. Sulzmann. An implementation of subtyping among regular expression types. In *APLAS 2004*, volume 3302 of *LNCS*, pages 57–73. Springer, 2004.
- [24] K. Matsuda and M. Wang. Bidirectionalization for Free with Runtime Recording: Or, a Light-weight Approach to the View-update Problem. In *PPDP 2013*, pages 297–308. ACM, 2013.
- [25] H. Pacheco and A. Cunha. Multifocal: A strategic bidirectional transformation language for XML schemas. In *ICMT 2012*, volume 7307 of *LNCS*, pages 89–104. Springer, 2012.
- [26] H. Pacheco, A. Cunha, and Z. Hu. Delta lenses over inductive types. In *BX 2012*, volume 49 of *Electronic Comms. of the EASST*, 2012.
- [27] H. Pacheco, Z. Hu, and S. Fischer. Monadic combinators for “putback” style bidirectional programming. In *PEPM 2014*, pages 39–50. ACM, 2014.
- [28] H. Pacheco, T. Zan, and Z. Hu. BiFluX: A Bidirectional Functional Update Language for XML. Technical Report GRACE-TR 2014-02, GRACE Center, National Institute of Informatics, Aug. 2014.
- [29] V. Preoteasa. A relation between unambiguous regular expressions and abstract data types. *Fundamentae Informatica*, 40(1):53–77, 1999.
- [30] J. Robie, D. Chamberlin, M. Dyck, D. Florescu, J. Melton, and J. Siméon. Xquery update facility 1.0. W3C Recommendation. <http://www.w3.org/TR/xquery-update-10/>, March 2011.
- [31] H. Seidl. On the finite degree of ambiguity of finite tree automata. *Acta Informatica*, 26(6):527–542, 1989.
- [32] S. Vansummeren. Type inference for unique pattern matching. *ACM TOPLAS*, 28(3):389–428, 2006.