# BETWEEN INVISIBLE XML AND VISIBLE XML

## Nico Verwer

# ixml parsers parse *text*

```
grammar  →  ixml parser generator
                              ↘
            plain text ⇢ ixml parser ⇢ XML
```
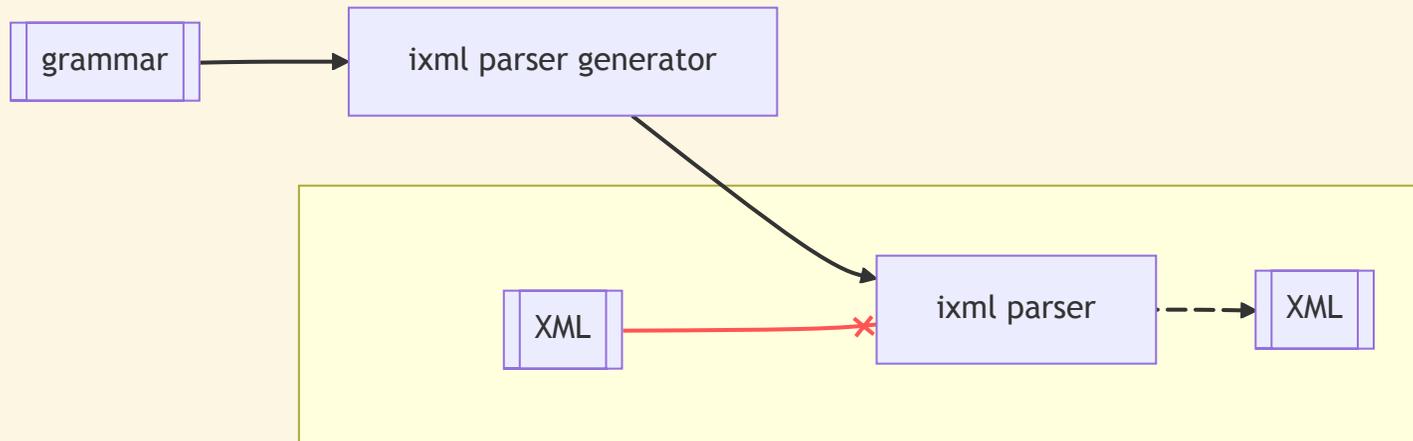
# real documents are XML ☺

Text fragments in the visible XML can contain invisible XML.

parsing → more structure → enriched documents

# what kind of 'real' documents?

> *" The 'Link eXtractor' (LX) recognizes references to legal sources, like legislation, court decisions, parliamentary documents and official publications in textual documents [...] by detecting fixed strings and rule-based patterns. "*

— https://gitlab.com/koop/ld/lx/lx-core

# invisible XML

'hello world'

```
date: day, s, month, s, year.
day: digit;
     digit, digit.
month: "January"; "February"; "March"; "April"; "May";
       "June"; "July"; "August"; "September"; "October";
       "November"; "December".
year: digit, digit, digit, digit.
-s: ([Zs]; #9; #a; #d)+ .
-digit: "0"; "1"; "2"; "3"; "4"; "5"; "6"; "7"; "8"; "9".
```

```
1 xquery version "3.1" encoding "UTF-8";
2
3 import module namespace toixml="http://rakensi.com/toixml" at "toixml.xqm";
4
5 (: invisible XML :)
6
7 let $grammar.ixml := ``[
8    date: day, s, month, s, year.
9    day: digit;
10        digit, digit.
11   month: "January"; "February"; "March"; "April"; "May"; "June";
12         "July"; "August"; "September"; "October"; "November"; "December".
13   year: digit, digit, digit, digit.
14   -s: ([Zs]; #9; #a; #d)+ .
15   -digit: "0"; "1"; "2"; "3"; "4"; "5"; "6"; "7"; "8"; "9".
16 ]``
17
18 let $parser as function(item()) as node()* := toixml:ixml-parser($grammar.ixml, ())
19
20 let $input as xs:string := '27 February 2026'
21
22 return $parser($input)
23
```

OK                                                                    5 : 2(

⊡  ⌂  Q   1 Result, 72 b                                              Result

```
<date>
  <day>27</day> <month>February</month> <year>2026</year>
</date>
```

# transparent XML

Suppose we have already recognized months.

```
<r>27 <month nr="2">February</month> 2026</r>
```

Keep the markup by making it *transparent.*

```
<r><date>
  <day>27</day>
  <month><month nr="2">February</month></month>
  <year>2026</year>
</date></r>
```

two <month> elements ☹

parse <month> again ☹

```
1 xquery version "3.1" encoding "UTF-8";
2
3 import module namespace toixml="http://rakensi.com/toixml" at "toixml.xqm";
4
5 (: transparent invisible XML :)
6
7 let $grammar.ixml := ``[
8   date: day, s, month, s, year.
9   day: digit;
10       digit, digit.
11   month: "January"; "February"; "March"; "April"; "May"; "June";
12          "July"; "August"; "September"; "October"; "November"; "December".
13   year: digit, digit, digit, digit.
14   -s: ([Zs]; #9; #a; #d)+ .
15   -digit: "0"; "1"; "2"; "3"; "4"; "5"; "6"; "7"; "8"; "9".
16 ]``
17
18 let $parser as function(item()) as node()* := toixml:ixml-parser($grammar.ixml, ())
19
20 let $input as element() := <r>27 <month nr="2">February</month> 2026</r>
21
22 return $parser($input)
23
```

OK                                                                          20 : 4

---

📥  ⌂  Q   1 Result, 121 b                                                  Result

```
<r>
  <date>
    <day>27</day> <month>
      <month nr="2">February</month>
    </month> <year>2026</year>
  </date>
</r>
```

# opaque XML

Use the `<month>` element that is already there by making it *opaque.*

tokenization, multi-stage parsing pipelines

```
date: day, s, <month>, s, year.
```

`<month>` is a 'pre-parsed non-terminal'

The parser checks if a `<month>` element is present at the right position.

One day, we might also match attributes.

```
1 xquery version "3.1" encoding "UTF-8";
2
3 import module namespace toixml="http://rakensi.com/toixml" at "toixml.xqm";
4
5 (: opaque invisible XML :)
6
7 let $grammar.ixml := ``[
8    date: day, s, <month>, s, year.
9    day: digit;
10       digit, digit.
11
12
13   year: digit, digit, digit, digit.
14   -s: ([Zs]; #9; #a; #d)+ .
15   -digit: "0"; "1"; "2"; "3"; "4"; "5"; "6"; "7"; "8"; "9".
16 ]``
17
18 let $parser as function(item()) as node()* := toixml:ixml-parser($grammar.ixml, ())
19
20 let $input as element() := <r>27 <month nr="2">February</month> 2026</r>
21
22 return $parser($input)
23
```

OK                                                                              8:3

---

⬇ ⌂ Q    1 Result, 94 b                                                    Result

```
<r>
  <date>
    <day>27</day> <month nr="2">February</month> <year>2026</year>
  </date>
</r>
```

# parsing *within* elements

```
<r>The symposium dates are
  <Date>26 <month nr="2">February</month> 2026</Date>
  and
  <Date>27 <month nr="2">February</month> 2026</Date>
</r>
```

```
let $options as map(*) := map {'parse-within-element' : 'Date'}
```

```
<r>The symposium dates are <Date>
    <date>
      <day>26</day> <month nr="2">February</month> <year>2026</ye
    </date>
  </Date> and <Date>
    <date>
      <day>27</day> <month nr="2">February</month> <year>2026</ye
    </date>
  </Date>
</r>
```

```
1 xquery version "3.1" encoding "UTF-8";
2
3 import module namespace toixml="http://rakensi.com/toixml" at "toixml.xqm";
4
5 (: opaque invisible XML within an element :)
6
7 let $grammar.ixml := ``[
8   date: day, s, <month>, s, year.
9   day: digit;
10        digit, digit.
11  year: digit, digit, digit, digit.
12  -s: ([Zs]; #9; #a; #d)+ .
13  -digit: "0"; "1"; "2"; "3"; "4"; "5"; "6"; "7"; "8"; "9".
14 ]``
15
16 let $options as map(*) := map {'parse-within-element' : 'Date'}
17 let $parser as function(item()) as node()* := toixml:ixml-parser($grammar.ixml, $options)
18
19 let $input as element() := <r>The symposium dates are <Date>26 <month nr="2">February</month> 2026</Date> and <
   Date>27 <month nr="2">February</month> 2026</Date></r>
20
21 return $parser($input)
22
```

OK                                                                                          19:1

---

🔽  ⌂  🔍  1 Result, 253 b                                                        Result

```
<r>The symposium dates are <Date>
    <date>
      <day>26</day> <month nr="2">February</month> <year>2026</year>
    </date>
  </Date> and <Date>
    <date>
      <day>27</day> <month nr="2">February</month> <year>2026</year>
    </date>
  </Date>
</r>
```

# invisible XML *matching*

```
<r>The symposium dates are
   26 <month nr="2">February</month> 2026 and
   27 <month nr="2">February</month> 2026
</r>
```

```
let $options as map(*) := map { 'complete-match' : 'false' }
```

```
<r>The symposium dates are <date>
     <day>26</day> <month nr="2">February</month> <year>2026</year
   </date> and <date>
     <day>27</day> <month nr="2">February</month> <year>2026</year
   </date>
</r>
```

This will find the longest matching fragments.

```
1 xquery version "3.1" encoding "UTF-8";
2
3 import module namespace toixml="http://rakensi.com/toixml" at "toixml.xqm";
4
5 (: invisible XML matching :)
6
7 let $grammar.ixml := ``[
8   date: day, s, <month>, s, year.
9   day: digit;
10       digit, digit.
11  year: digit, digit, digit, digit.
12  -s: ([Zs]; #9; #a; #d)+ .
13  -digit: "0"; "1"; "2"; "3"; "4"; "5"; "6"; "7"; "8"; "9".
14 ]``
15
16 let $options as map(*) := map { 'complete-match' : 'false' }
17 let $parser as function(item()) as node()* := toixml:ixml-parser($grammar.ixml, $options)
18
19 let $input as element() := <r>The symposium dates are 26 <month nr="2">February</month> 2026 and 27 <month nr="2">
   February</month> 2026</r>
20
21 return $parser($input)
22
```

◎ OK                                                                                                        19 : 1

⎙  ⌂  Ｑ   1 Result, 203 b                                                                                   Result

```
<r>The symposium dates are <date>
    <day>26</day> <month nr="2">February</month> <year>2026</year>
  </date> and <date>
    <day>27</day> <month nr="2">February</month> <year>2026</year>
  </date>
</r>
```
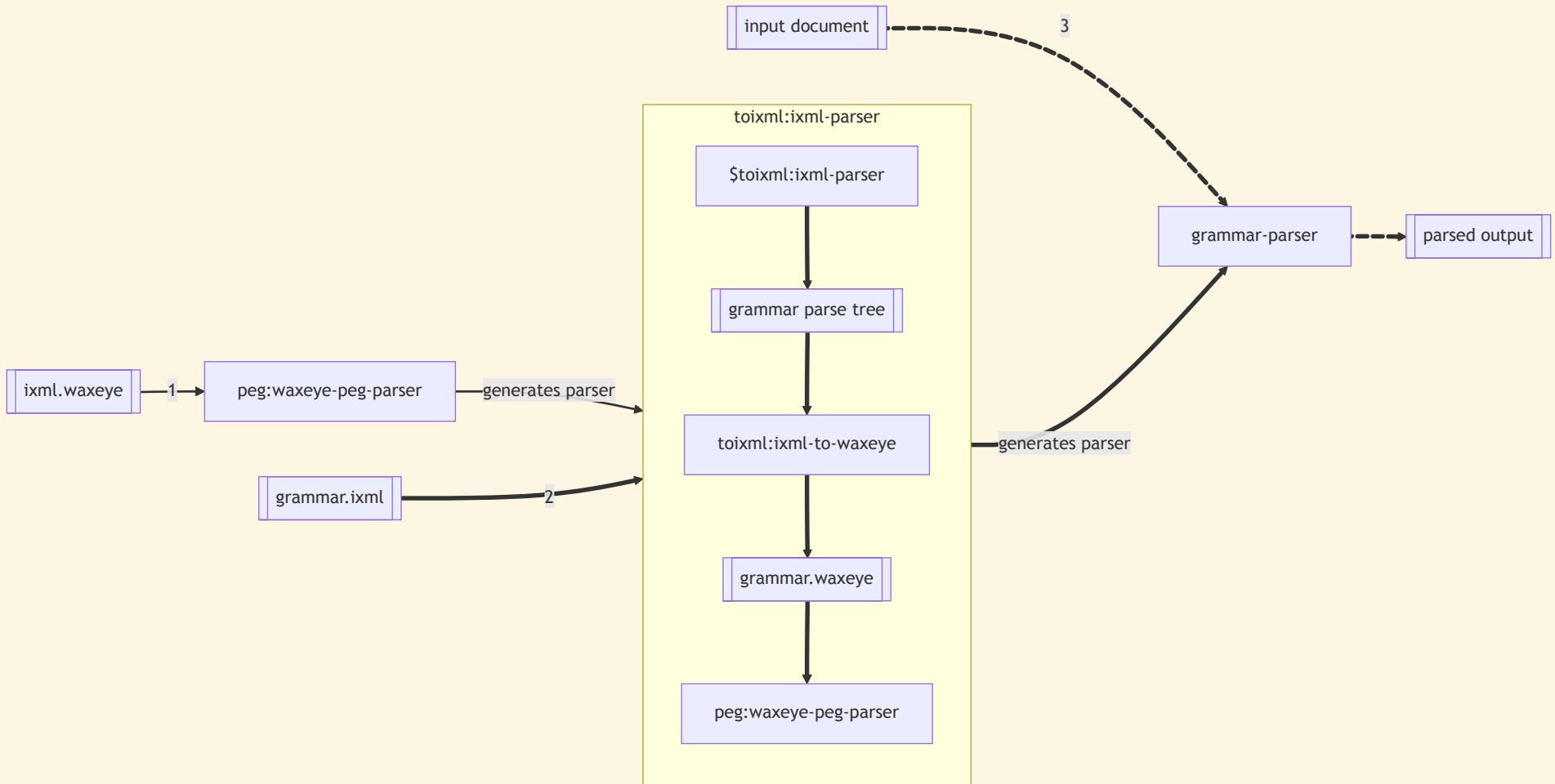
# implementation

LX uses a PEG parser generator (waxeye), not ixml.

Working on an ixml - PEG implementation,
which will be open source.

The waxeye implementation is not open source.

The examples use a *transpiler*.

not compliant (@, -, +, character classes)

```
                                    input document ┄┄┄┄┄┄┄┄┄┄ 3
                                                              ┄┄┄┄┄┄┄┄┄┄┐
                    ┌─────────────────────────────────────────┐         ▼
                    │         toixml:ixml-parser               │   ┌──────────────┐      ┌──────────────┐
                    │                                          │   │grammar-parser│ ┄┄┄> │parsed output │
                    │      ┌──────────────────────┐            │   └──────────────┘      └──────────────┘
                    │      │  $toixml:ixml-parser │            │         ▲
                    │      └──────────────────────┘            │         │
                    │                 │                        │         │
                    │                 ▼                        │         │
                    │      ┌──────────────────────┐            │         │
                    │      │  grammar parse tree  │            │         │
                    │      └──────────────────────┘            │         │
                    │                 │                        │         │
                    │                 ▼                        │         │
ixml.waxeye ─1─> peg:waxeye-peg-parser  generates parser ─────>│ toixml:ixml-to-waxeye │──generates parser─┘
                    │      └──────────────────────┘            │
                    │                 │                        │
                    │                 ▼                        │
grammar.ixml ───2────>  │      ┌──────────────────────┐            │
                    │      │   grammar.waxeye     │            │
                    │      └──────────────────────┘            │
                    │                 │                        │
                    │                 ▼                        │
                    │      ┌──────────────────────┐            │
                    │      │ peg:waxeye-peg-parser │            │
                    │      └──────────────────────┘            │
                    └─────────────────────────────────────────┘
```

# if there is time left

negation (Fredrik Öhrström's talk)

partial binding / parsing pipeline (John Lumley's talk)

```xquery
1  xquery version "3.1" encoding "UTF-8";
2
3  import module namespace toixml="http://rakensi.com/toixml" at "toixml.xqm";
4
5  (: negation (negative look-ahead), inspired by Fredrik Öhrström's talk :)
6  (: This is cheating, because a PEG parser would choose the longest match for word and number anyway. :)
7
8  let $grammar.ixml := ``[
9    target : (word ; number)+ .
10   word : letter+, !letter .
11   number : digit+, !digit .
12   -letter : ["a"-"z"] .
13   -digit: ["0"-"9"] .
14 ]``
15
16 let $parser as function(item()) as node()* := toixml:ixml-parser($grammar.ixml, ())
17 let $input as xs:string := 'alfa123beta456'
18 return $parser($input)
19
```

OK                                                                                                    10 : 23

---

🔖  △  🔍  1 Result, 104 b                                                                            Result

```xml
<target>
  <word>alfa</word>
  <number>123</number>
  <word>beta</word>
  <number>456</number>
</target>
```

```
39    quantity: digits .
40    -digits: ["0"-"9"]+ .
41    -s: -" "+.
42    -lf: -#d ; -#a .
43 ]``
44
45 let $customer-parser as function(item()) as node()* :=
46   toixml:ixml-parser($customer-grammar.ixml, map { 'complete-match' : 'false' })
47
48 let $order-parser as function(item()) as node()* :=
49   toixml:ixml-parser($order-grammar.ixml, ())
50
51 let $input as element() :=
52 <orders>
53 Steven Pemberton
54 21 Sandridge Road
55 St Albans AL1 4BY
56 United Kingdom
57 item:ballpoint pen,3
58 item:notebook,1
59 </orders>
60
61 return $input => $customer-parser() => $order-parser()
62
```

OK                                                                    5:71

⊡  ⌂  Q   1 Result, 502 b                                        Result

```
<orders>
  <order>
<customer>
      <person>
        <given>Steven</given> <surname>Pemberton</surname>
      </person>
<street>
        <no>21</no> <streetname>Sandridge Road</streetname>
      </street>
```

# LX's parsing pipeline (small fragment)

```
( $input

=> (: named entity recognition for legislation / regulations :)
  transform:if( $extract-leg-reg,
    lx-ner:regeling(?, $is-test, $is-recognize)
  )

=> (: parse references to EU legislation / regulations :)
  transform:if( $extract-eu-leg-reg,
    ( lx-peg:eu-regelgeving(?)
    , transform:xsl('links/xsl/normalize-eu-regelgeving.xslt', ())
    , lido:resolve-celex(?, $options?test-url)
    )
  )

=> (: local alias detection :)
  transform:if( $extract-leg-reg,
    ( lx-peg:lokale-alias(?)
```

# conclusion

invisible XML: parse text input into XML output

visible XML: parse text from XML input

transparent XML: ignore and keep XML structure

opaque XML: pre-parsed non-terminals

parse within specific elements

match invisible XML fragments