# Design Suggestions for XML++

XML Prague 14.02.2020

# Ruthless Dictator Slogan

# Make XML Great Again!

# XML Design flaws

**DoS Attack**

Expand to >3 gigabytes memory

```
1  <?xml version="1.0"?>
2  <!DOCTYPE lolz [
3   <!ENTITY lol "lol">
4   <!ELEMENT lolz (#PCDATA)>
5   <!ENTITY lol1 "&lol;&lol;&lol;&lol;&lol;&lol;&lol;&lol;&lol;&lol;">
6   <!ENTITY lol2 "&lol1;&lol1;&lol1;&lol1;&lol1;&lol1;&lol1;&lol1;&lol1;&lol1;">
7   <!ENTITY lol3 "&lol2;&lol2;&lol2;&lol2;&lol2;&lol2;&lol2;&lol2;&lol2;&lol2;">
8   <!ENTITY lol4 "&lol3;&lol3;&lol3;&lol3;&lol3;&lol3;&lol3;&lol3;&lol3;&lol3;">
9   <!ENTITY lol5 "&lol4;&lol4;&lol4;&lol4;&lol4;&lol4;&lol4;&lol4;&lol4;&lol4;">
10  <!ENTITY lol6 "&lol5;&lol5;&lol5;&lol5;&lol5;&lol5;&lol5;&lol5;&lol5;&lol5;">
11  <!ENTITY lol7 "&lol6;&lol6;&lol6;&lol6;&lol6;&lol6;&lol6;&lol6;&lol6;&lol6;">
12  <!ENTITY lol8 "&lol7;&lol7;&lol7;&lol7;&lol7;&lol7;&lol7;&lol7;&lol7;&lol7;">
13  <!ENTITY lol9 "&lol8;&lol8;&lol8;&lol8;&lol8;&lol8;&lol8;&lol8;&lol8;&lol8;">
14  ]>
15  <lolz>&lol9;</lolz>
```

See Wiki: One Billion Laugh Attack / XML Bomb

# Cut old Braids



German Expression: **Cut old braids** - Students during Wars of Liberation in 1813-1815

# Simplify XML

- MicroXML (2012) solves XML Security & Complexity:
  - **MicroXML prohibits 12 features of XML 5th edition**
    - No custom entities
    - Only UTF-8
    - …
    - No draconian XML error handling (see XML 5 - 2014-2016)
  - See James Clark Blog Entry (Dec. 2010)
  - Better marketing: LeanXML, SecureXML,…, XML2.0
  - Desired: Transformations from old to new XML

- **Summary:** Fix/Simplify XML for new implementations!

## Simplify Tooling

- # Why do we need **>3** XML grammars?
  - **RelaxNG** (rng) used by ODF
  - **W3C Schema** (xsd) used by
    - OOXML
    - UBL
    - UN/CEFACT Cross-Industry Invoice (CII) XML
  - **DTD** explicitly not used by HTML 5

# Valid HTML 5

```
1  <!DOCTYPE html>
2  <html lang="en">
3      <head style="color:red">
4          <meta charset="UTF-8"/>
5          <title>No grammar restrictions @HTML</title>
6      </head>
7      <body style="color:red">
8          <p>Hello World!</p>
9      </body>
10 </html>
```

- DTD unset as not sufficient enough
- HTML 5 allows styles everywhere

# The Power of XML grammar

- Most expressive: RelaxNG

- Less expressive: W3C Schema

- Least expressive: DTD (should be dropped)

See [Taxonomy of XML Schema Languages Using Formal Language Theory](#) (2005)
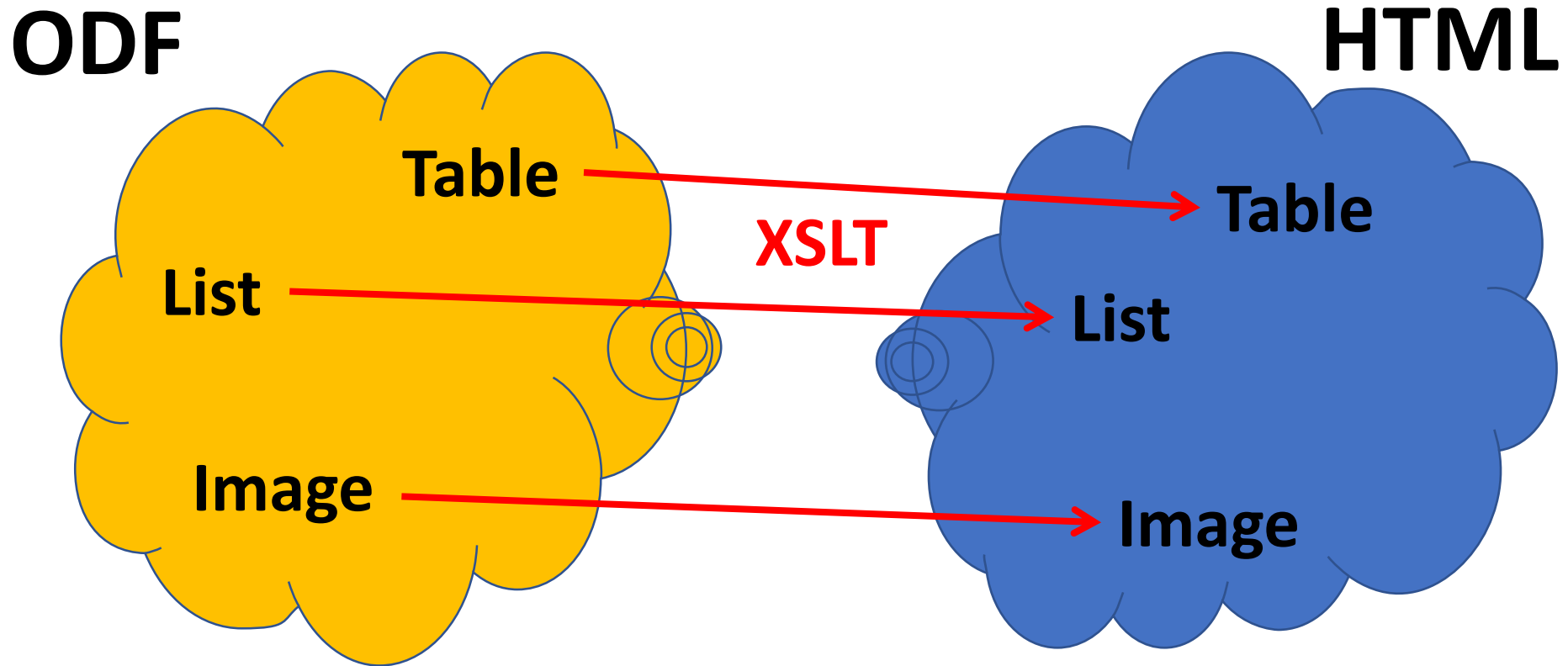
# XML Grammar Tools

- Sun's Multi-Schema-Validator ([MSV](#))

  - Internal RelaxNG Model (most expressive)
  - Loads DTD, W3C Schema & RelaxNG

- [Trang Converter](#) from James Clark

  - W3C XML Schema only as output! (why?)

- How to generate scanners/validators from grammar:
  [Regular-Expression Derivatives](#) (paper from 1964)

**Journey into my Past**

# 5 User Requirements

# Example Requirements (1/5) - Mapping
**Sun's Web Office (1999) – ODF Viewer**

- **Web Viewer of StarOffice (later ODF) Documents**

**ODF**

**HTML**

**Table** → **Table**

**XSLT**

**List** → **List**

**Image** → **Image**

# Example Requirements (1/5) - Mapping

**Sun's Web Office (1999) – ODF Viewer**

- Management & myself were often asking:

  - How much is done? When are we ready? (Grammar can tell coverage! – Saxon feature?)
  - Where is the bug? In my mapping or my XSLT?
  - Reuse some work for a Java DOM / SAX approach?

**Web Office 2nd Edition (till 2011)**
**Elaborate a mapping based on the Grammar**

- **Web Editor ODF Documents – both way mapping**

**ODF**

**HTML**

**Table**

**Table**

**Java**

**List**

**List**

**Image**

**Image**

# Example Requirements (3/5) – XML Run Time Model

**Run Time Model from RelaxNG Grammar**

- Java XML Binding worked only for W3C Schema

- Load, Edit and Save ODF

- **ODF Runtime Environment –** ODFDOM (Java)

  - Loaded RelaxNG into <u>MultiSchemaValidator (MSV)</u>
  - Using MSV API to fill <u>templates by (Apache Velocity)</u>
  - Created <u>from grammar Java (typed) DOM classes</u>: *<draw:image>* → *DrawImageElement class*

# Example Requirements (3/5) – XML Run Time Model
**Run Time Model from RelaxNG Grammar**

- XML DOM not succinct data structure

- Not well for huge documents

- Not streaming XML by Netflix

- Spreadsheet Pixel-Art problematic:

# Example Requirements (4/5) – RTM close Semantic

**Syntax is an implementation detail, Semantic counts!**

- I've been so proud of ODF XML..

- [J David Eisenberg](#) once said in a conf call:
*"Leave me alone with your XML details, I want easy*
***semantics****: Open Text doc, add paragraph with 'HelloWorld!'*
*No implementation details desired!"*

- ***JAXB*** *created RunTimeModels overtaking XML syntax as*
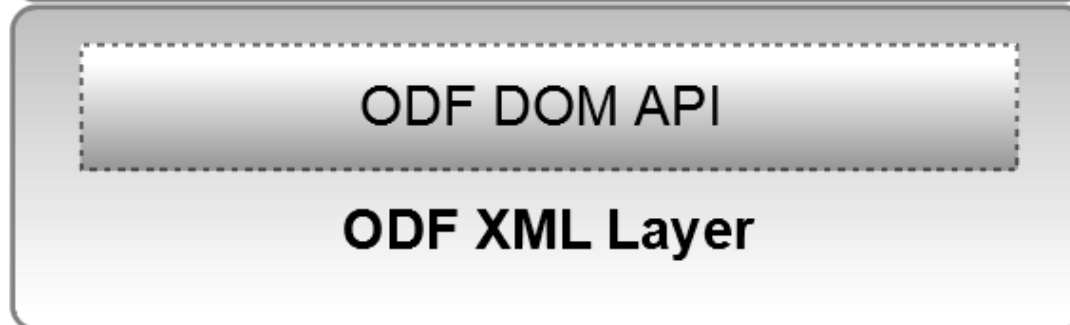*semantic (improvable by abstracting XML details)*

Svante Schubert

# Example Requirements (5/5) – Generate Everything

**<u>Generate as much as possible!</u>**



**3.** ODF User API / ODF Semantic Layer

**2.** ODF DOM API / ODF XML Layer
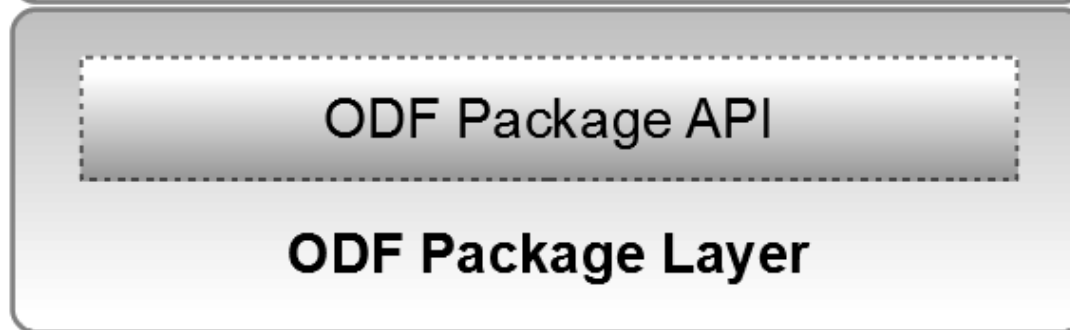
**1.** ODF Package API / ODF Package Layer

Svante Schubert

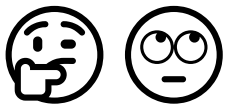# Example Requirements (5/5) – Generate Everything

**Generation is based on Grammar**

- ODF RelaxNG grammar is a text file of > 18.000 lines

- Hard to extract information by ODF users

- Infos missing: Semantic User Entities consisting of multiple XML like a table not yet defined in grammar!

# ODF GRAMMAR
**HARD TO ANSWER**

Can a

paragraph **<text:p>**
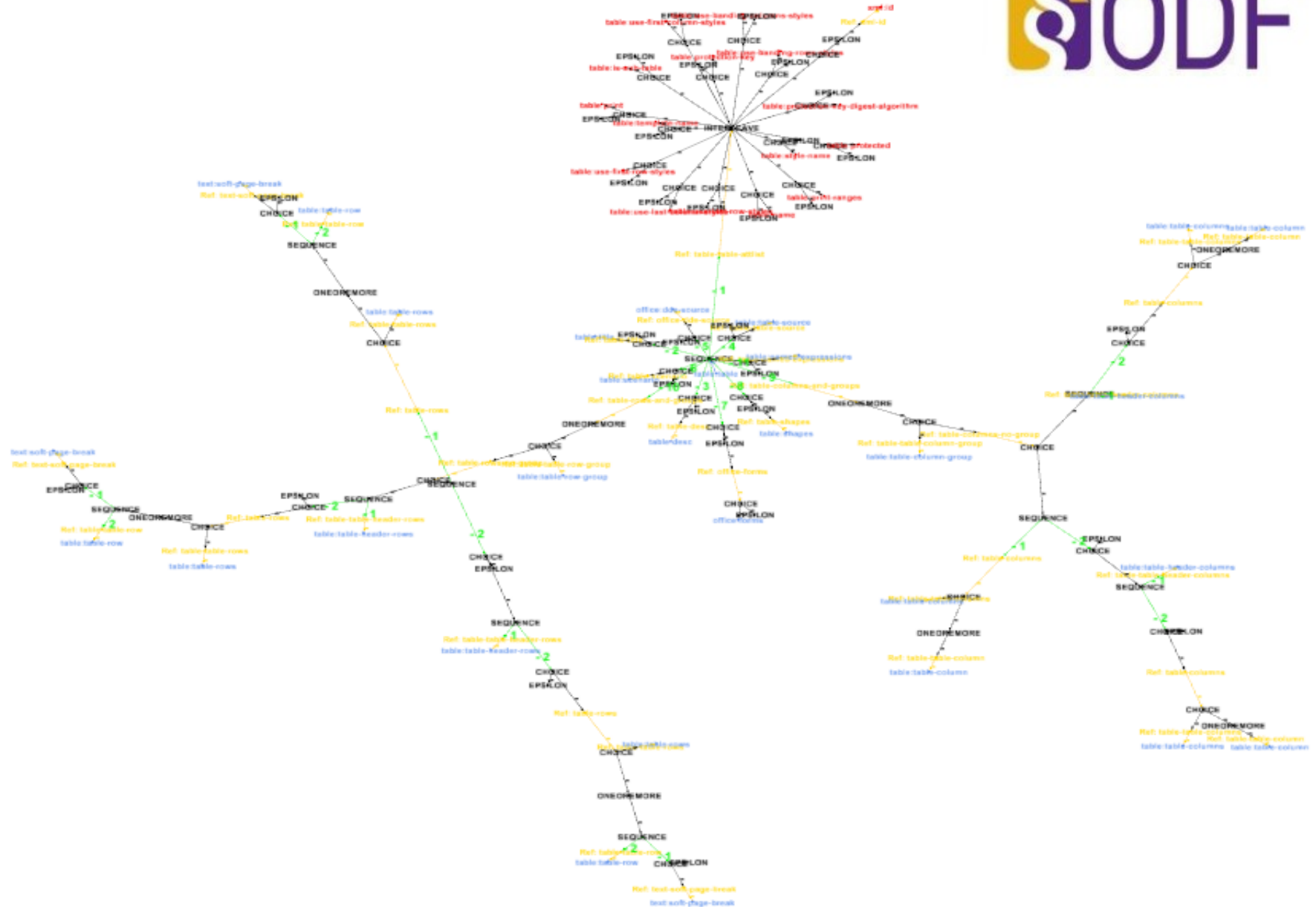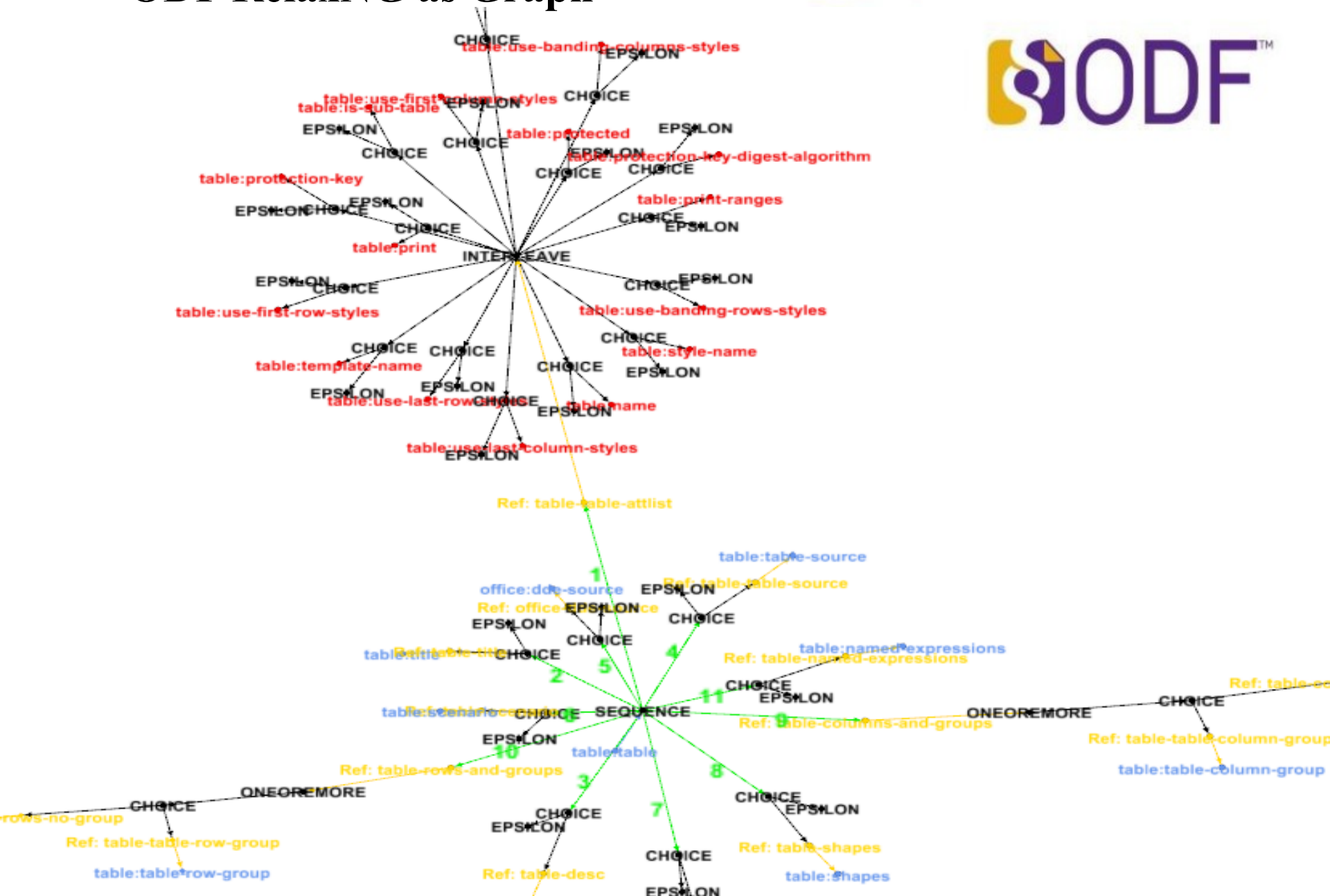be nested
in a valid document?
😗 🙄

**ODF 1.2 XML**:

- **598 XML Elements**

- **1300 XML Attributes**

>18k lines

# Example Requirements (5/5) – Generate Everything
**Generation based as RelaxNG Graph**

- Solvable by traversing a (tree-ish) graph…

- Based on idea: Loading source code for analysis as graphs

- Loaded ODF RelaxNG from MSV into Apache Tinkerpop

- Saved Subgraphs from element to element as GraphML

- The *<table:table>* element with all its child elements

- Table GraphML rendered in Gephi

# ODF RelaxNG as Graph

# ODF RelaxNG as Graph

# ODF RelaxNG as Graph

# ODF RelaxNG as Graph

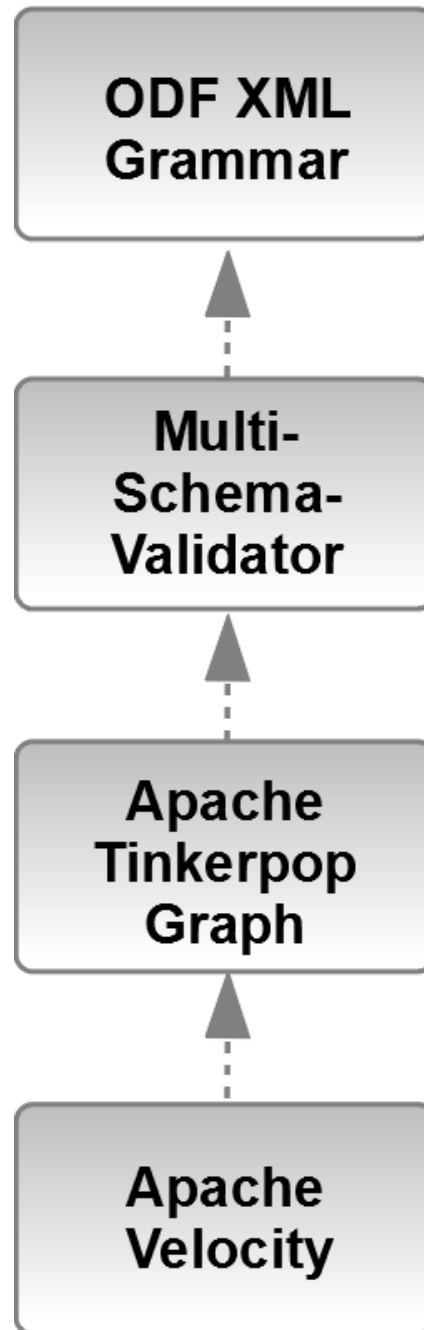text:table-row

text:soft-page-break

**1**

**2**

**SEQUENCE**

# Example Requirements (5/5) –Generate Everything
## Improved RunTimeModel from XML Grammar

- JSON is far easier usable as RunTimeModel than XML

- Some Grammar to RunTimeModel features are still missing:

  - Insertion handling for multiple optional child elements: A,B,C
  - Parent with multiple similar children with ID
    Parent should become a lazy map (created on demand)
  - General missing XML grammar annotation/reverse-engineering Tool

Svante Schubert

# ODF Toolkit

**Source Code**

**Generator**

**Architecture**

```
┌──────────────┐
│   ODF XML    │
│   Grammar    │
└──────────────┘
        ▲
        ┊
┌──────────────┐
│    Multi-    │
│   Schema-    │
│  Validator   │
└──────────────┘
        ▲
        ┊
┌──────────────┐
│    Apache    │
│  Tinkerpop   │
│    Graph     │
└──────────────┘
        ▲
        ┊
┌──────────────┐
│    Apache    │
│   Velocity   │
└──────────────┘
```
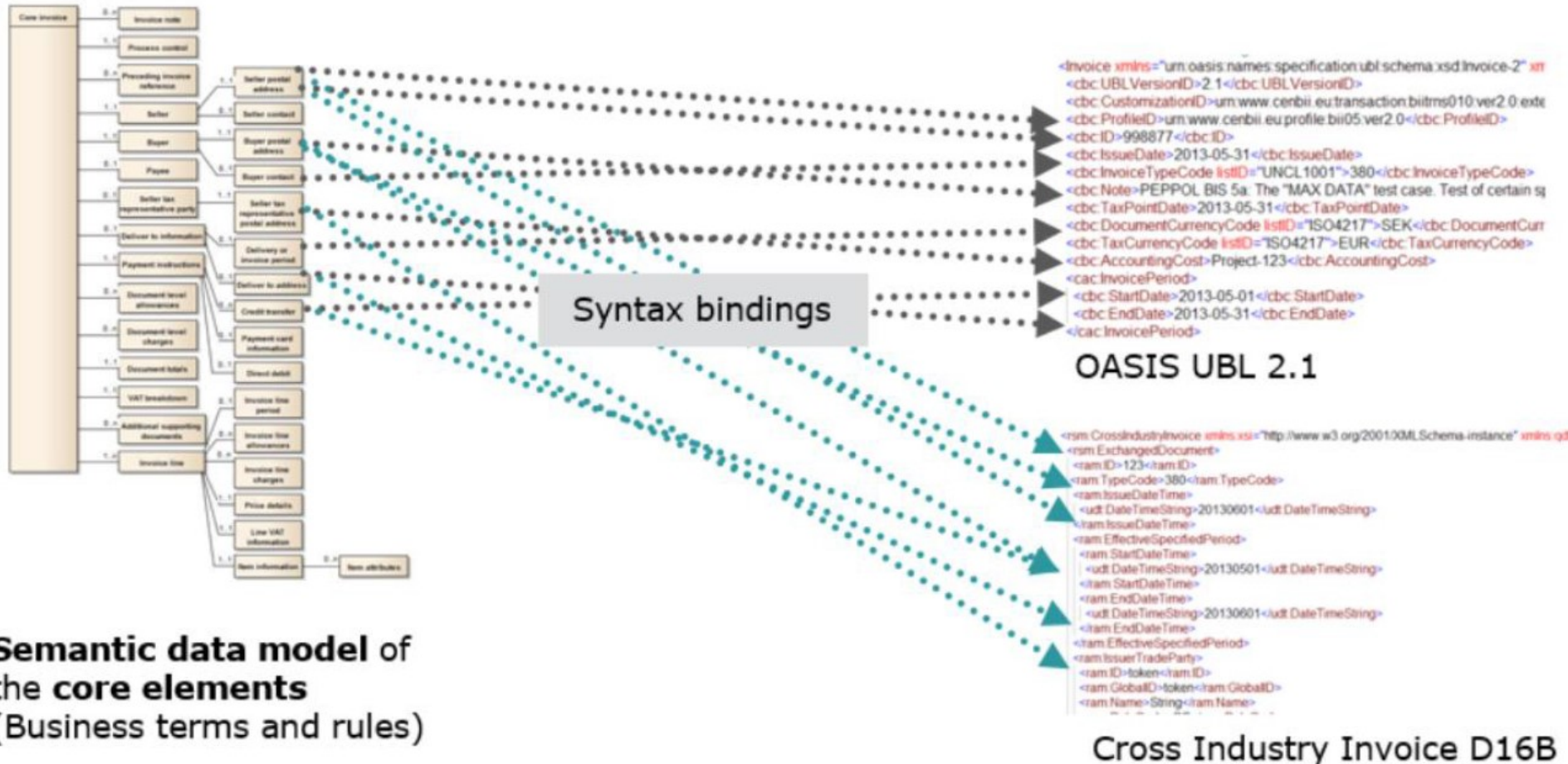
**RelaxNG grammar**:
~18000 text lines
~600 XML elements
~1200 XML attributes

XML Valdiator -
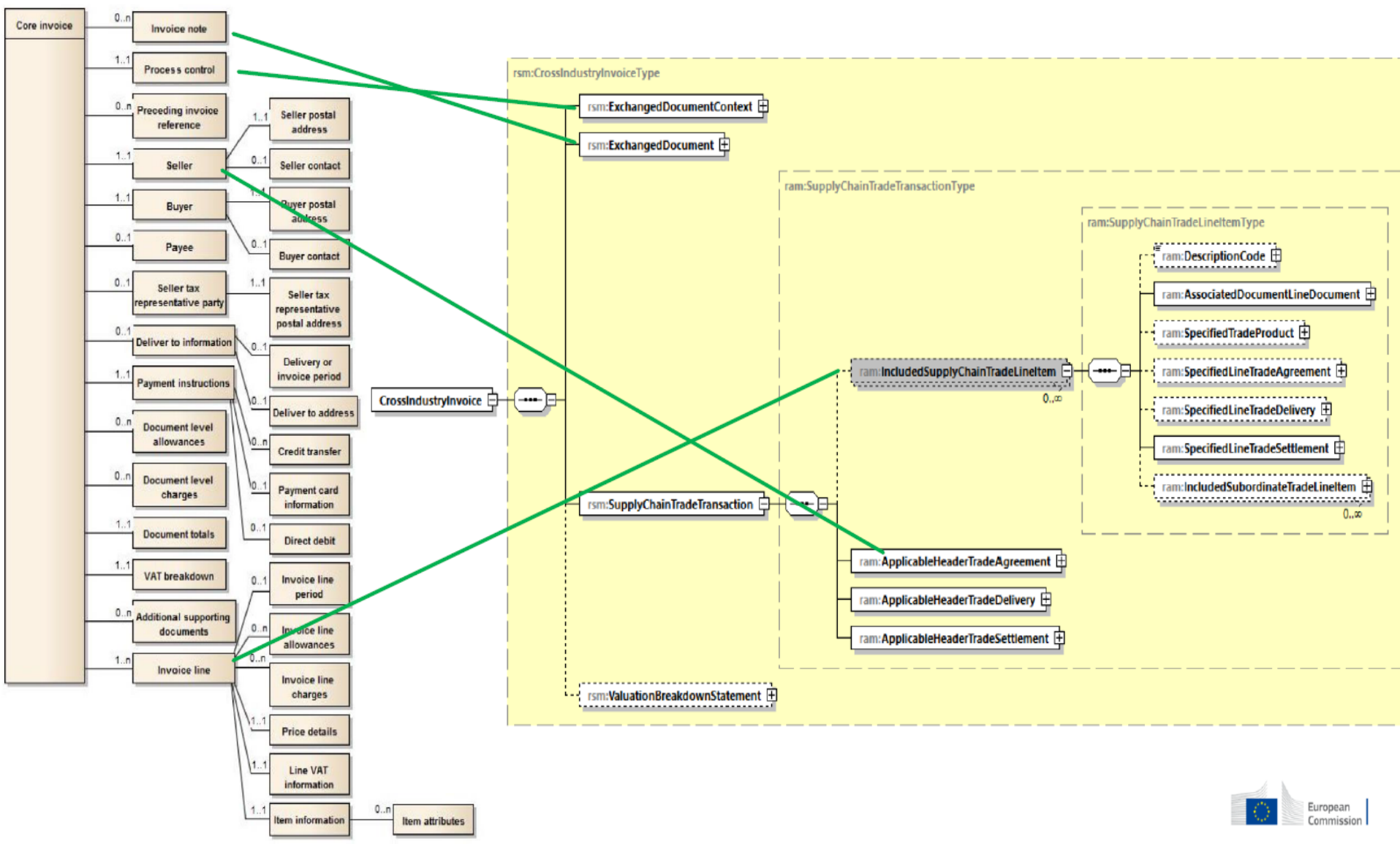**reads** many
**XML grammar**s

**Graph** of
XML grammar
**(since ODF TOOLKIT 1.0.0)**

**template engine**
- generating sources
by text templates with context
with Java access
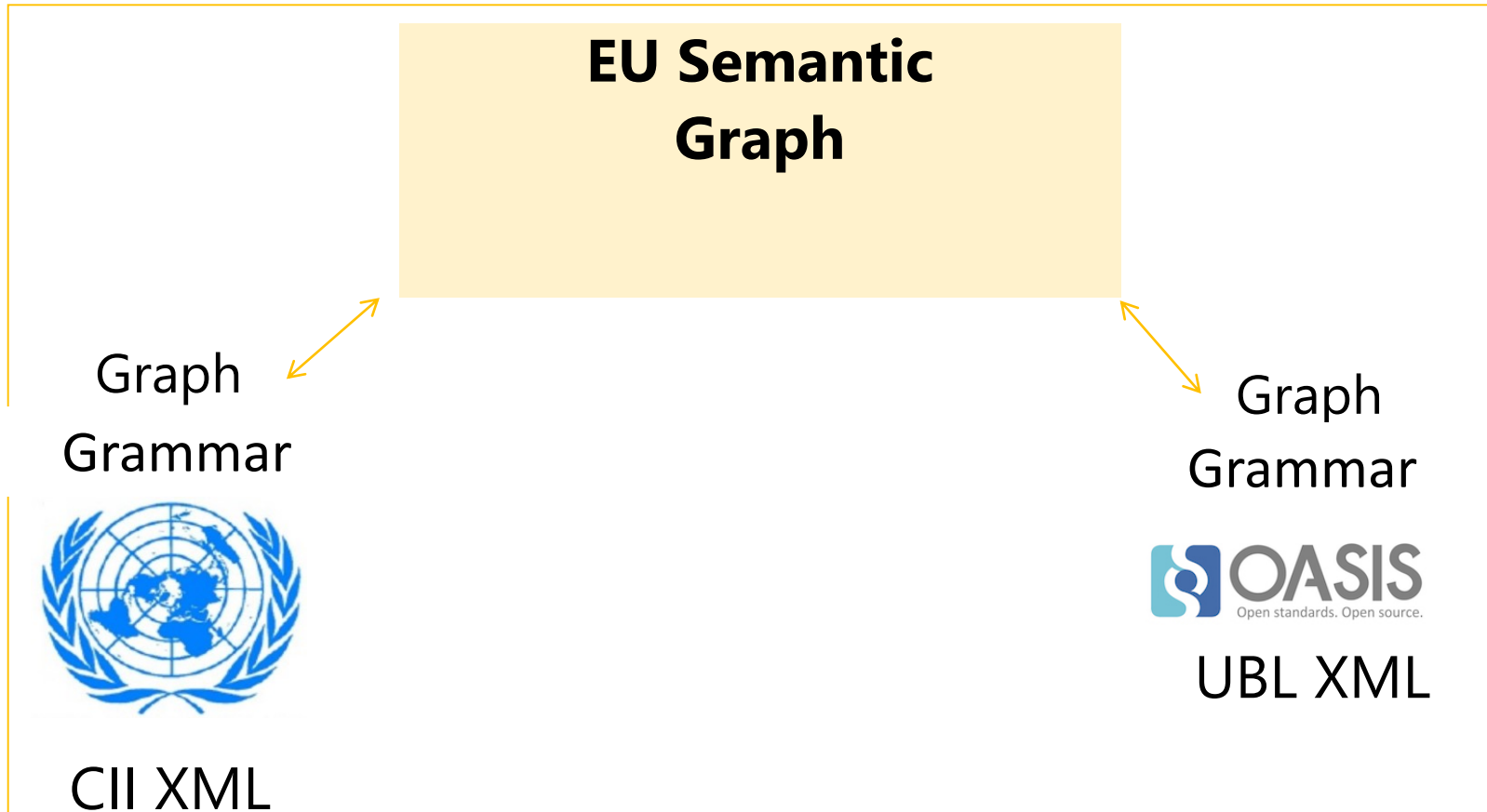
# Semantic Model
# as Glue between XML Formats



Semantic data model of
the core elements
(Business terms and rules)

Syntax bindings

OASIS UBL 2.1

Cross Industry Invoice D16B

# Semantic Model to CII XML Format

# Graph as Data Model

Data from EU Specification as Graph
(Sum of 3 connected Graphs)

**EU Semantic Graph**

Graph Grammar

Graph Grammar

UBL XML

CII XML

## Take Away

- XML should be simplified for newcomers

- XML Grammar are far more powerful than expected

- Too little Grammar conversion & tooling

- **Focus** should be more **on Semantic** instead of Syntax


To be continued…