

# From monolithic XML for print/web to lean XML for data: realising linked data for dictionaries

**Matt Kohl & Sandro Cirulli**  
**Language Technologists**  
**Oxford University Press (OUP)**

7 June 2014







- ▶ **Print-oriented:** designed to capture dictionary layout
- ▶ **Monolithic:** one enormous document
- ▶ **Permissive:** continually loosened to accommodate new texts

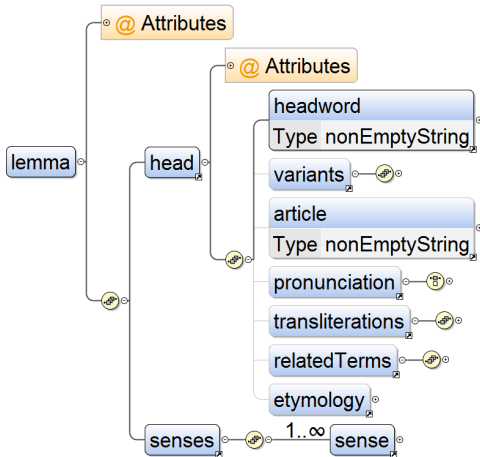
*Can't give us the flexibility we need*

A new approach should:

- ▶ **Represent language concepts**, not layouts
- ▶ **Enable data reusability** for different products & services
- ▶ **Allow only one, clear way to model any given lexical item**

# Data Modelling

## The New Lexical Schema



### Conversion Framework Requirements

- ▶ **Scalability:** convert 40+ data-sets
- ▶ **Standardization:** harmonize variation inside the data-sets
- ▶ **Modularity:** enable customization, slotting in & out of QA, etc.

# Data Conversion

## Tools

---

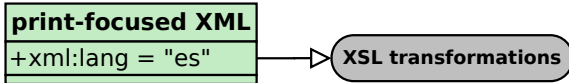
- ▶ XProc
- ▶ XSpec
- ▶ Schematron & XML Schema
- ▶ Jenkins CI
- ▶ Agile methodology





# Data Conversion

## Simplified XProc pipeline



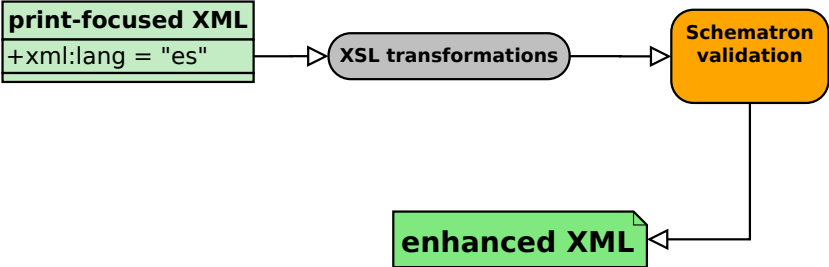
# Data Conversion

## Simplified XProc pipeline



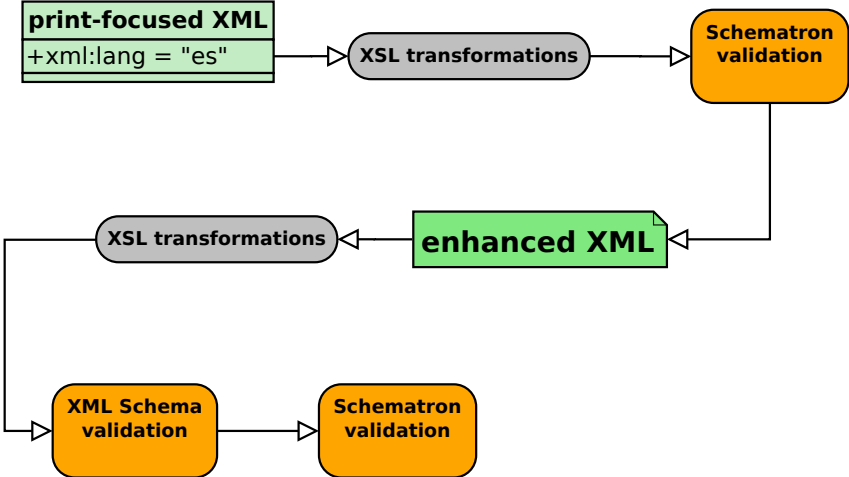
# Data Conversion

## Simplified XProc pipeline



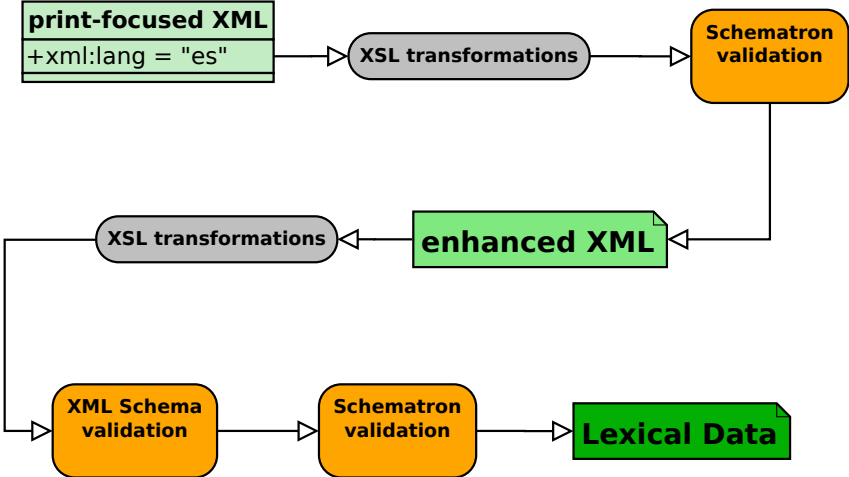
# Data Conversion

## Simplified XProc pipeline



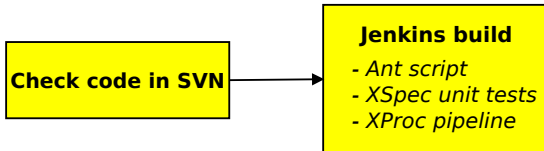
# Data Conversion

## Simplified XProc pipeline



# Data Conversion

## Build Workflow



# Data Conversion

## Build Workflow

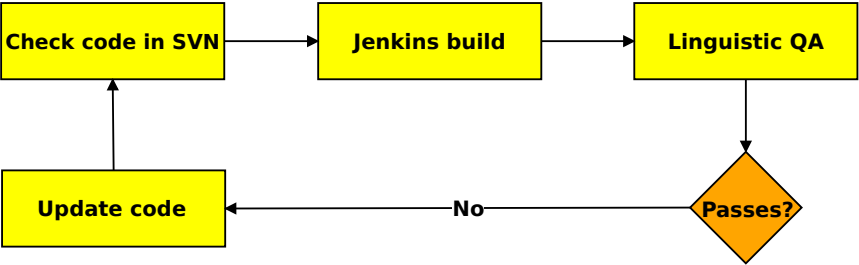
---





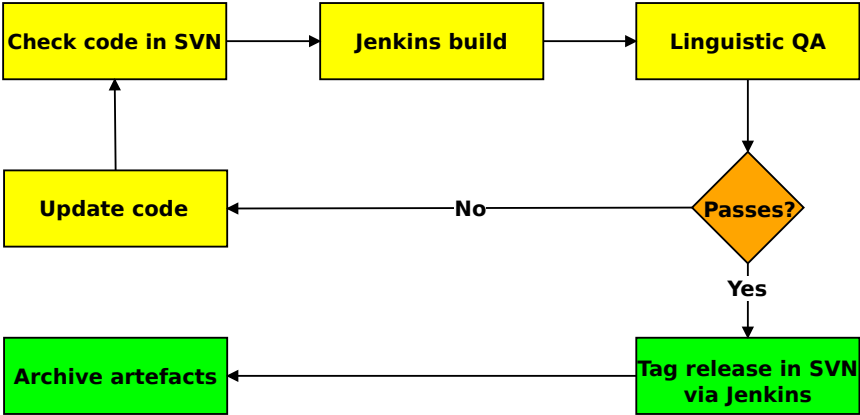
# Data Conversion

## Build Workflow



# Data Conversion

## Build Workflow





### Print-focused DTD

```
<se2 num="2">  
  <lg>  
    <ge>Esp</ge>  
    <reg>coloquial</reg>  
  </lg>  
  <msDict type="core">  
    <df>Persona que  
      tiene mal  
      carácter o mala  
      intención.</df>  
    <syn>malaleche</syn>  
  </msDict>  
</se2>
```

### Print-focused DTD

```
<se2 num="2">
  <lg>
    <ge>Esp</ge>
    <reg>colloquial</reg>
  </lg>
  <msDict type="core">
    <df>Persona que
      tiene mal
      carácter o mala
      intención.</df>
    <syn>malaleche</syn>
  </msDict>
</se2>
```

### New Lexical XSD

```
<sense register="informal"
  region="ES">
  <definitions>
    <definition>
      <text>Persona que tiene
        mal carácter o mala
        intención</text>
    </definition>
  </definitions>
  <synonyms>
    <synonym>malaleche</
      synonym>
  </synonyms>
</sense>
```

### Print-focused DTD

```
<se2 num="2">
  <lg>
    <ge>Esp</ge>
    <reg>coloquial</
      reg>
  </lg>
  <msDict type="core">
    <df>Persona que
      tiene mal
      carácter o mala
      intención.</df>
    <syn>malaleche</syn>
  </msDict>
</se2>
```

### New Lexical XSD

```
<sense
  register="informal"
  region="ES">
  <definitions>
    <definition>
      <text>Persona que tiene
        mal carácter o mala
        intención</text>
    </definition>
  </definitions>
  <synonyms>
    <synonym>malaleche</
      synonym>
  </synonyms>
</sense>
```



# Next Steps

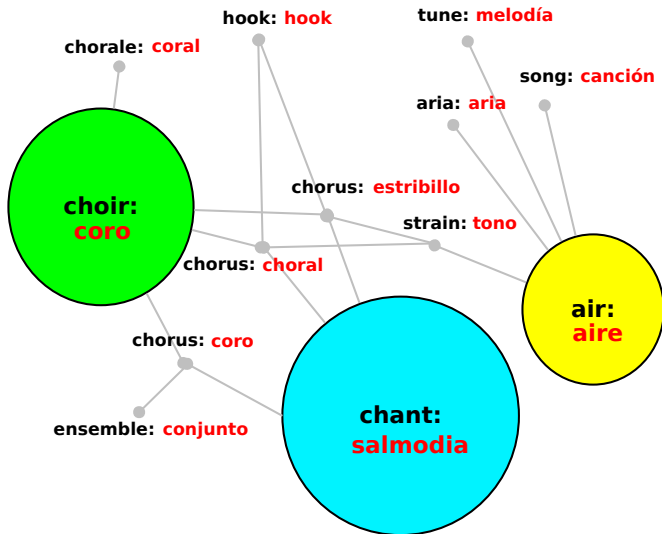
## Prototype RDF/XML

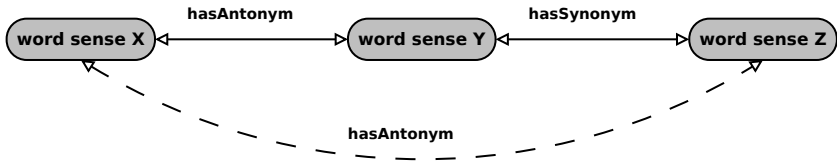
```
<Sense rdf:about="sense:es_noun_malauva_se_2">
  <isDescribedBy
    rdf:resource="
      definition:es_noun_malauva_se_2_def_1"/>
  <hasRegister rdf:resource="register:informal"
    />
  <hasRegion rdf:resource="region:ES"/>
  <hasSynonym rdf:resource="lemma:a5e644"/>
</Sense>
<StandardDefinition
  rdf:about="definition:es_noun_malauva_se_2_def_1">
  <rdfs:label xml:lang="es">Persona que tiene mal
    carácter o mala intención</rdfs:label>
</StandardDefinition>
```



# RDF Data extraction

## Musical terms in English & Spanish





## ▶ Overall project requirements

- ▶ Moving from products to platforms and services
- ▶ Supporting current business needs while innovating
- ▶ Adapting in nimble ways to fast changing market requirements
- ▶ Focusing on time and cost efficiency

## ▶ Data model

- ▶ Content driven
- ▶ Machine interpretable
- ▶ Modular
- ▶ Evolvable/adaptable

## ▶ Conversion process

- ▶ Highly automated
- ▶ Modular
- ▶ Scalable

**Thank you for your attention!**  
**Any questions?**

**Matt Kohl: [matt.kohl@oup.com](mailto:matt.kohl@oup.com)**  
**Sandro Cirulli: [sandro.cirulli@oup.com](mailto:sandro.cirulli@oup.com)**

