

# **Inferring XML Schema Definitions from XML Data**

Geert Jan Bex, Frank Neven and  
Stijn Vansummeren

Hasselt University and transnational  
University of Limburg, Belgium

# Overview

- **Introduction**
- Complete algorithm *iLOCAL*
- Heuristic *iXSD*
- Experiments
- Conclusions

# Motivation for schemas

- Why schemas?
  - automation & optimization of search
  - integration of XML data sources
  - translation & processing of XML data
  - used by software tools, e.g., JAXB, Castor
  - schema matching & model management
- Why infer schemas?
  - 50 % of XML document on the web have none [Barbosa et al., 2005]
  - 33 % of schemas are not valid [Bex et al., 2004, 2005]

real world XML & XSDs

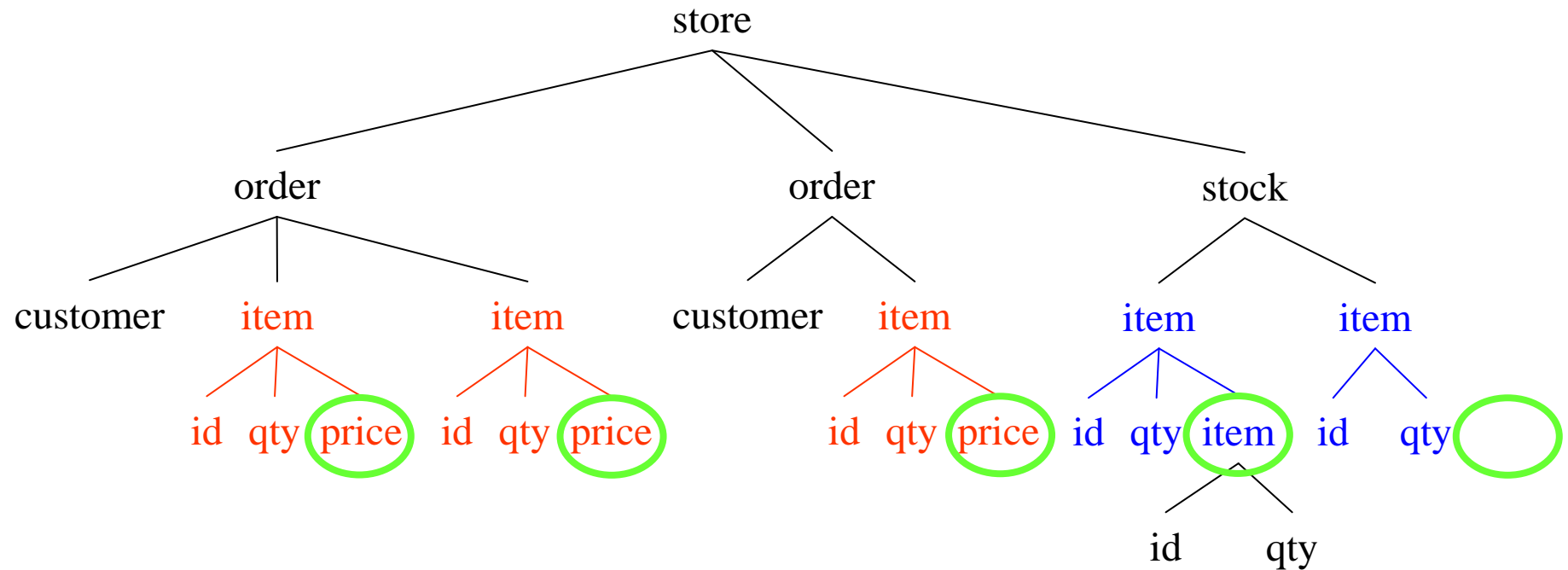
# Motivation for XSD inference

- DTD inference
  - XTract [Garofalakis et al., 2003]
  - trang [Clark]
  - *i*DTD [Bex et al., 2006]
- XSD inference
  - trang
  - XStruct
  - JAXB, .Net

} output XSD syntax,  
but equivalent to DTD

expressive power limited to that of DTDs!

# How do DTDs and XSDs differ?



in DTDs, either:

item → id, qty, (price + item\*)

or

order\_item → id, qty, price

stock\_item → id, qty, stock\_item\*

can be done in XSDs

# XSD: abstract syntax

```
<xsd:element name="store" type="store"/>
```

```
<xsd:complexType name="store">
```

```
  <xsd:sequence>
```

```
    <xsd:element name="order" type="order" minOccurs="0" maxOccurs="unbounded"/>
```

```
    <xsd:element name="stock" type="stock"/>
```

```
  </xsd:sequence>
```

```
</xsd:complexType>
```

```
<xsd:complexType name="order">
```

```
  <xsd:sequence>
```

```
    <xsd:element name="customer" type="customer"/>
```

```
    <xsd:element name="item" type="item1" minOccurs="1" maxOccurs="unbounded"/>
```

```
  </xsd:sequence>
```

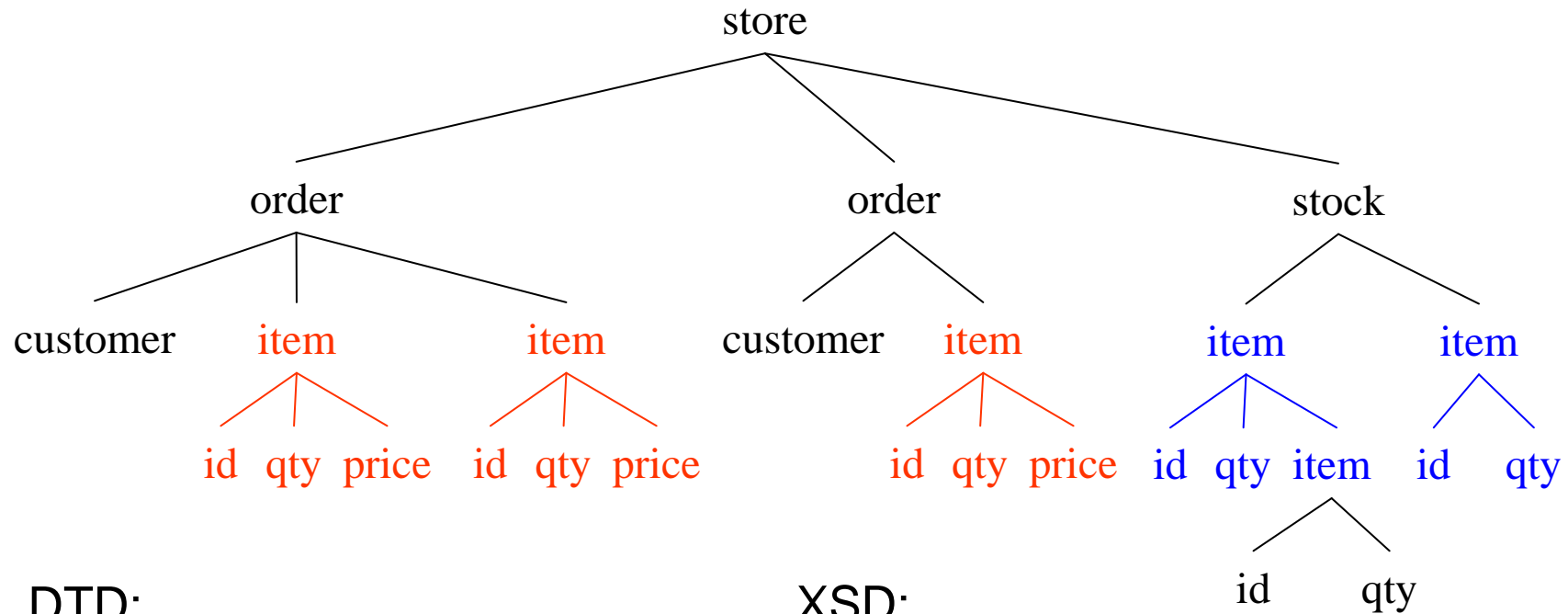
```
</xsd:complexType>
```

```
root → store[store]
```

```
store → order[order]*, stock[stock]
```

```
order → customer[customer], item[item1]+
```

# Motivating example for XSD



DTD:

**root** → **store**

**store** → **order\***, **stock**

**order** → **customer**, **item<sup>+</sup>**

**stock** → **item<sup>+</sup>**

**item** → **id**, **qty**, (**price** + **item\***)

XSD:

*root* → **store**[*store*]

*store* → **order**[*order*]\*, **stock**[*stock*]

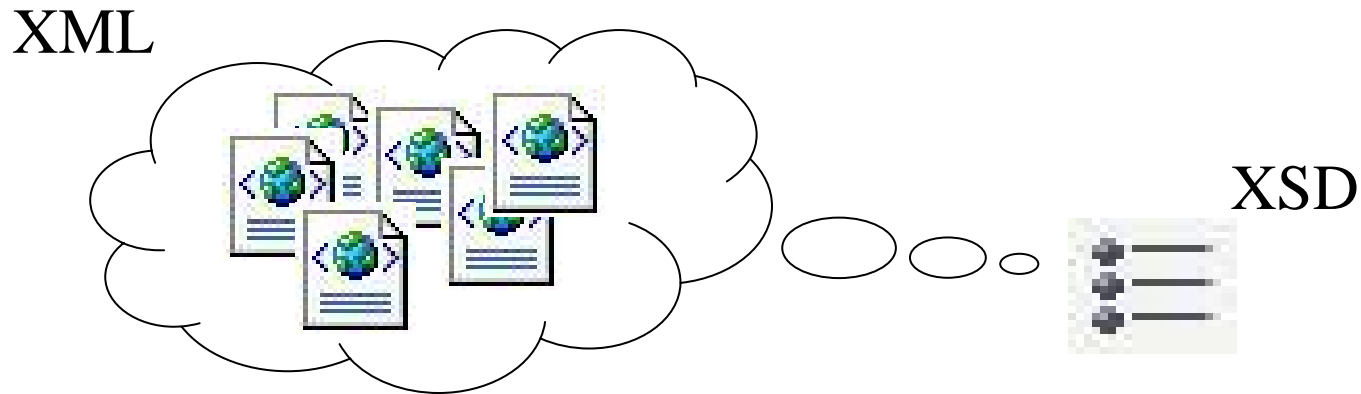
*order* → **customer**[*customer*], **item**[*item*<sub>1</sub>]<sup>+</sup>

*stock* → **item**[*item*<sub>2</sub>]<sup>+</sup>

*item*<sub>1</sub> → **id**[*id*], **qty**[*qty*], **price**[*price*]

*item*<sub>2</sub> → **id**[*id*], **qty**[*qty*], **item**[*item*<sub>2</sub>]\*

# Inference of XSDs



- Problem: infer XSD from XML corpus
- Requirement: concise, i.e., humans can interpret/validate
- But... theorem [Gold, 1967]:

impossible to learn from positive data only



# XSD property

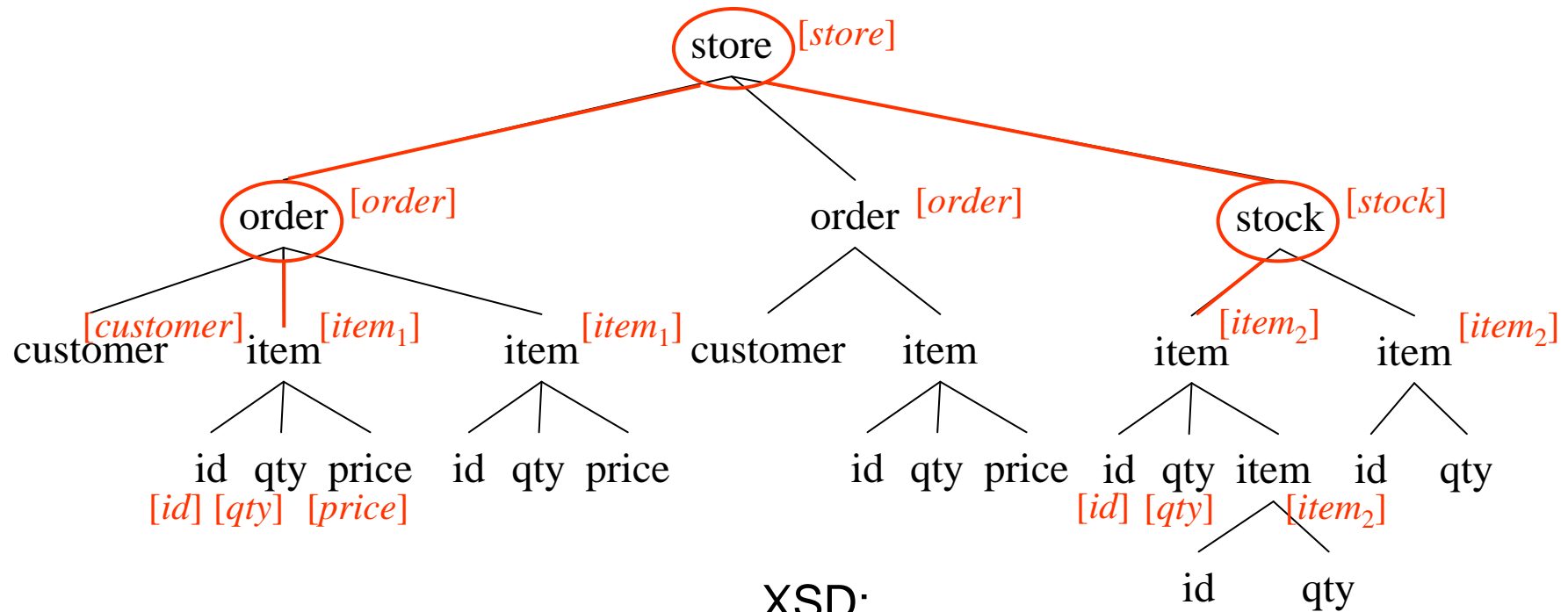
W3C specs: **Element Declarations Consistent (EDC)**:

no elements with distinct type in same content model

~~$some\ type \rightarrow item\ [item_1]^+, item\ [item_2]^+$~~

content model of an element depends on its context

# XML validation for XSD



XSD:

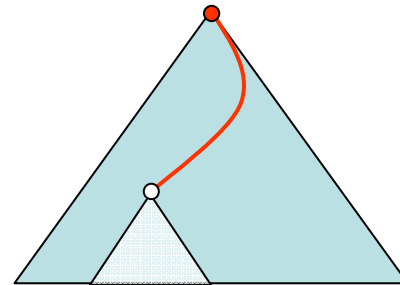
$root \rightarrow \mathbf{store}[store]$   
 $store \rightarrow \mathbf{order}[order]^*, \mathbf{stock}[stock]$   
 $order \rightarrow \mathbf{customer}[customer], \mathbf{item}[item_1]^+$   
 $stock \rightarrow \mathbf{item}[item_2]^+$   
 $item_1 \rightarrow \mathbf{id}[id], \mathbf{qty}[qty], \mathbf{price}[price]$   
 $item_2 \rightarrow \mathbf{id}[id], \mathbf{qty}[qty], \mathbf{item}[item_2]^*$

if XML is valid:  
type assignment is determined  
by path from element to root

# XML validation for XSD

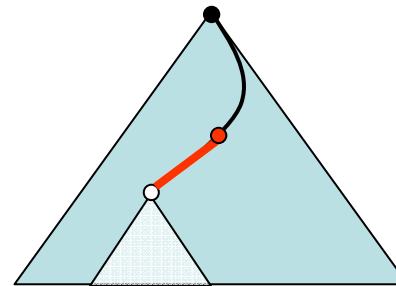
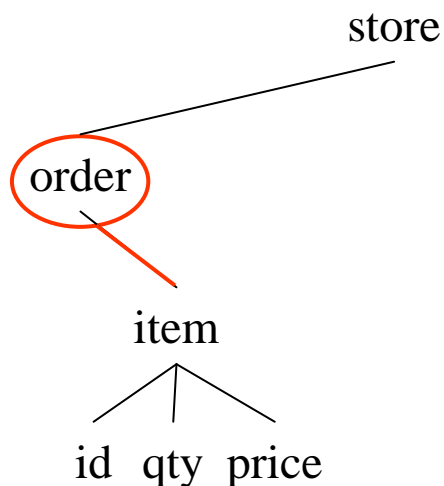
Theorem [Martens et al., 2006]

Content model of an element is uniquely determined by the path from the root to that element



# XSD observations: local context

- Large, diverse corpus of real world XSDs [Bex et al., 2004, Martens et al., 2006]
  - **98 %** of XSDs only local context: relevant ancestor path has length of **at most 3**, i.e., "greatgrandfather"



# XSD observations: SOREs

- Large, diverse corpus of real world XSDs [Bex et al., 2004, Martens et al., 2006]
  - 99 % of regular expressions is single occurrence

- What's a Single Occurrence RegExp

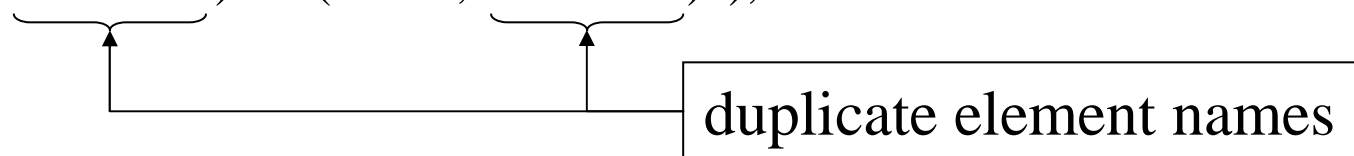
header, protein, organism, reference\*, comment\*, genetics\*, complex\*, function\*, classification?, keywords?, feature\*, summary, sequence

authors, citation, volume?, month?, year, pages?, (title + descr)?, xrefs?

title, (author, affiliation?)<sup>+</sup>, abstract

- ... and what's not

title, ((author, affiliation)<sup>+</sup> + (editor, affiliation)<sup>+</sup>), abstract



# Overview

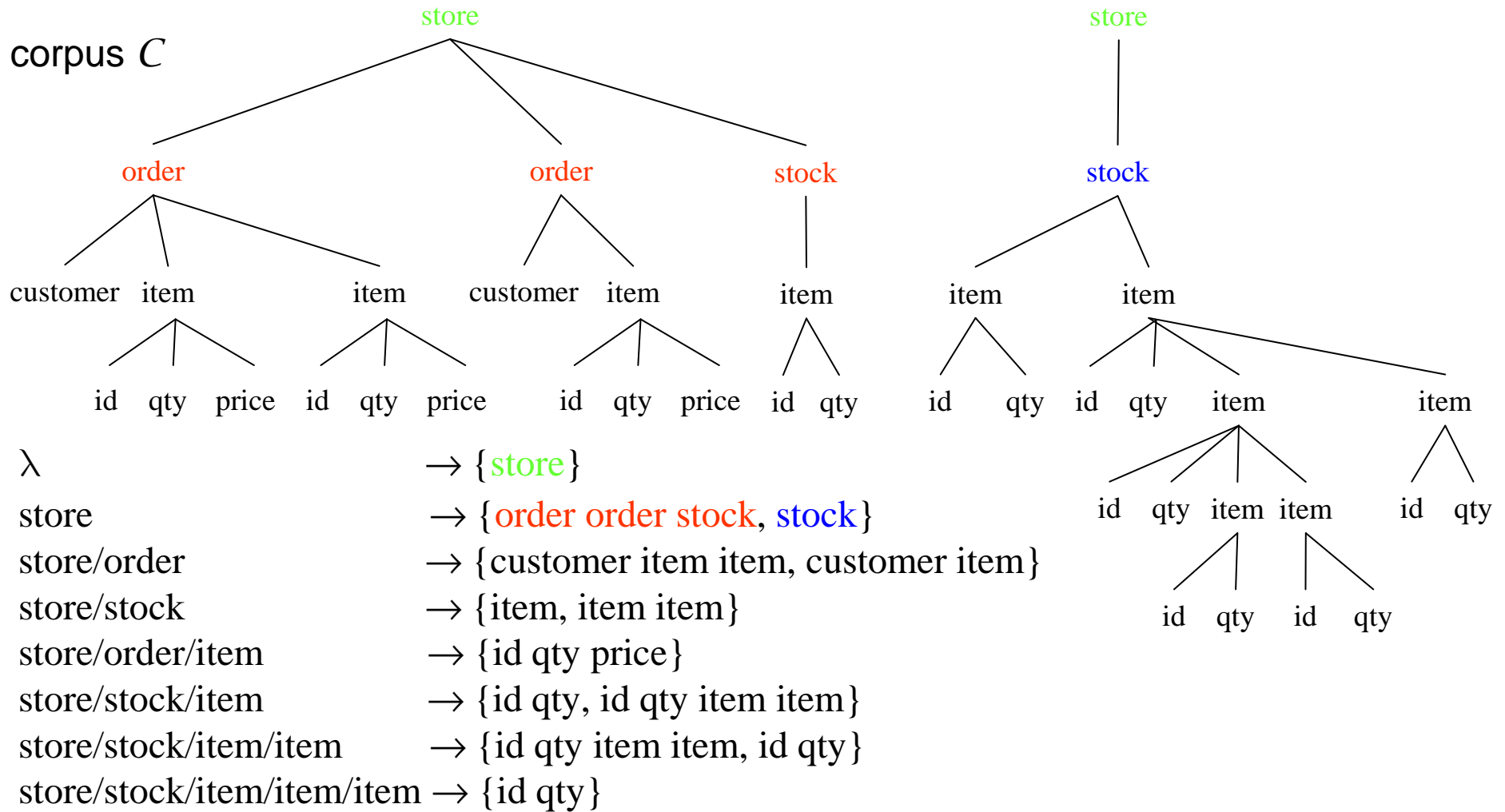
- Introduction
- Complete algorithm *iLOCAL*
- Heuristic *iXSD*
- Experiments
- Conclusions

# Main result

Theorem:

XSDs with local context and SORE content models  
are learnable from positive examples only

# Algorithm *i*LOCAL



paths are types [Martens et al., 2006]



# Algorithm *i*LOCAL

$\lambda$	$\rightarrow$ {store}
store	$\rightarrow$ {order order stock, stock}
store/order	$\rightarrow$ {customer item item, customer item}
store/stock	$\rightarrow$ {item, item item}
store/order/item	$\rightarrow$ {id qty price}
store/stock/item	$\rightarrow$ {id qty, id qty item item}
store/stock/item/item	$\rightarrow$ {id qty item item, id qty}
store/stock/item/item/item	$\rightarrow$ {id qty}

locality:  $k = 2$



$\lambda$	$\rightarrow$ {store}
store	$\rightarrow$ {order order stock, stock}
store/order	$\rightarrow$ {customer item item, customer item}
store/stock	$\rightarrow$ {item, item item}
order/item	$\rightarrow$ {id qty price}
stock/item	$\rightarrow$ {id qty, id qty item item}
item/item	$\rightarrow$ {id qty item item, id qty}

# Algorithm *i*Local

$\lambda$	→ {store}
store	→ {order order stock, stock}
store/order	→ {customer item item, customer item}
store/stock	→ {item, item item}
order/item	→ {id qty price}
stock/item	→ { <b>id qty, id qty item item</b> }
item/item	→ {id qty item item, id qty}



*i*SOA, ToSORE [Bex et al., 2006]

XSD

<i>root</i>	→ <b>store</b> [ <i>store</i> ]
<i>store</i>	→ <b>order</b> [ <i>store/order</i> ]*, <b>stock</b> [ <i>store/stock</i> ]
<i>store/order</i>	→ <b>customer</b> [ <i>order/customer</i> ], <b>item</b> [ <i>order/item</i> ] <sup>+</sup>
<i>store/stock</i>	→ <b>item</b> [ <i>stock/item</i> ] <sup>+</sup>
<i>order/item</i>	→ <b>id</b> [ <i>item/id</i> ], <b>qty</b> [ <i>item/qty</i> ], <b>price</b> [ <i>item/price</i> ]
<i>stock/item</i>	→ <b>id</b> [ <i>item/id</i> ], <b>qty</b> [ <i>item/qty</i> ], <b>item</b> [ <i>item/item</i> ]*
<i>item/item</i>	→ <b>id</b> [ <i>item/id</i> ], <b>qty</b> [ <i>item/qty</i> ], <b>item</b> [ <i>item/item</i> ]*

# Algorithm *i*LOCAL

- Theorem: *i*LOCAL is sound

corpus  $C$  is valid with respect to inferred XSD

- Theorem: *i*LOCAL is  $k$ -complete

if corpus  $C$  is "sufficiently large" then  
target XSD is equivalent with inferred XSD

# Algorithm MINIMIZE

<i>root</i>	→ <b>store</b> [ <i>store</i> ]	
<i>store</i>	→ <b>order</b> [ <i>store/order</i> ]*, <b>stock</b> [ <i>store/stock</i> ]	
<i>store/order</i>	→ <b>customer</b> [ <i>order/customer</i> ], <b>item</b> [ <i>order/item</i> ] <sup>+</sup>	
<i>store/stock</i>	→ <b>item</b> [ <i>stock/item</i> ] <sup>+</sup>	
<i>order/item</i>	→ <b>id</b> [ <i>item/id</i> ], <b>qty</b> [ <i>item/qty</i> ], <b>price</b> [ <i>item/price</i> ]	
<i>stock/item</i>	→ <b>id</b> [ <i>item/id</i> ], <b>qty</b> [ <i>item/qty</i> ], <b>item</b> [ <i>item/item</i> ]*	} duplicate types
<i>item/item</i>	→ <b>id</b> [ <i>item/id</i> ], <b>qty</b> [ <i>item/qty</i> ], <b>item</b> [ <i>item/item</i> ]*	

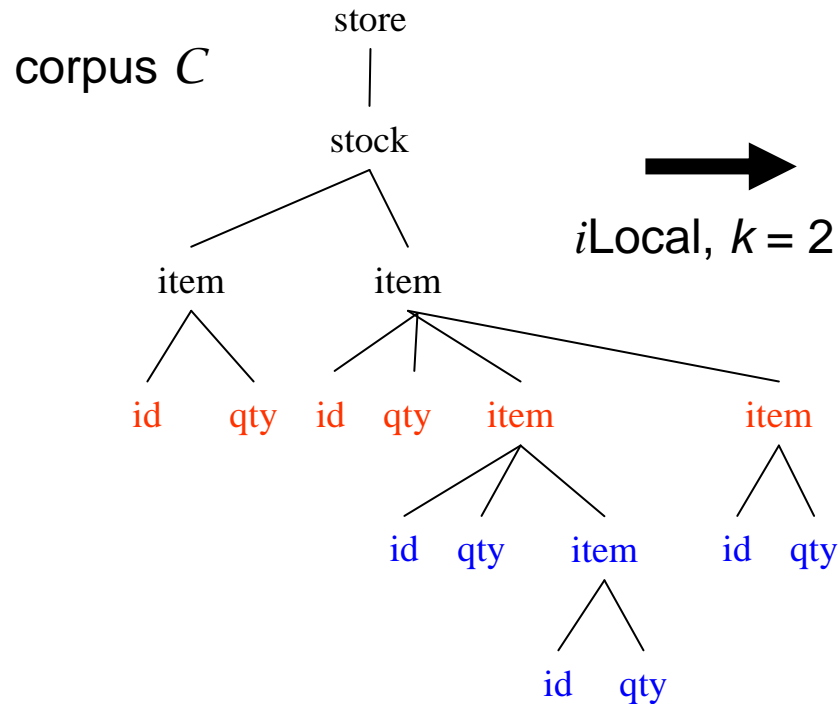
MINIMIZE ↓

<i>root</i>	→ <b>store</b> [ <i>store</i> ]
<i>store</i>	→ <b>order</b> [ <i>store/order</i> ]*, <b>stock</b> [ <i>store/stock</i> ]
<i>store/order</i>	→ <b>customer</b> [ <i>order/customer</i> ], <b>item</b> [ <i>order/item</i> ] <sup>+</sup>
<i>store/stock</i>	→ <b>item</b> [ <i>item<sub>2</sub></i> ] <sup>+</sup>
<i>order/item</i>	→ <b>id</b> [ <i>item/id</i> ], <b>qty</b> [ <i>item/qty</i> ], <b>price</b> [ <i>item/price</i> ]
<i>item<sub>2</sub></i>	→ <b>id</b> [ <i>item/id</i> ], <b>qty</b> [ <i>item/qty</i> ], <b>item</b> [ <i>item<sub>2</sub></i> ]*

# Overview

- Introduction
- Complete algorithm *i*LOCAL
- **Heuristic *i*XSD**
- Experiments
- Conclusions

# In practice: incomplete data



$stock/item \rightarrow \{id\ qty, id\ qty\ item\ item\}$   
 $item/item \rightarrow \{id\ qty\ item, id\ qty\}$

$iSOA, ToSORE \downarrow$

$stock/item \rightarrow id[item/id], qty[item/qty],$   
 $item[item/item]^*$   
 $item/item \rightarrow id[item/id], qty[item/qty],$   
 $item[item/item]?$

MINIMIZE can't minimize!

incomplete data  $\Rightarrow$  *iLocal* derives too many types!

# Practical heuristics

- Define "distance" between types
  - details: see paper
- For types  $s, t$ : if  $\text{distance}(s, t) < \varepsilon$ , unify  $s$  and  $t$   
= REDUCE
- Our practical algorithm  $iXSD$ :

$$iXSD(k, C) = \text{REDUCE}(iLOCAL(k, C))$$

# Overview

- Introduction
- Complete algorithm *i*LOCAL
- Heuristic *i*XSD
- **Experiments**
- Conclusions



# Experiments

- Corpora:

- 697 real world XSD documents:  $C_{\text{XSD}}$

real world corpora  
are hard to find

- XSD schema is local with  $k = 2$
    - **attributeGroup, group, extension**: 2 contexts
    - **restriction**: 3 contexts

- 8 corpora for synthetic XSDs, 200 XML documents each:  $C_1, \dots, C_8$

- XSD schemas define documents of unbounded depth, width
    - local with  $k = 2, 3$
    - 12 to 23 types
    - one schema associates multiple types with six element names
    - XML generated with ToXgene

# Precision

types of *i*XSD imprecisions:

1. **content model** for target and inferred type can differ
2. type in target XSD can correspond to multiple types in inferred XSD: **false positives**
3. type in inferred XSD can correspond to multiple types in target XSD: **false negatives**
4. type in target XSD is **not derived**  
*incomplete corpus, can't be avoided*

# 1. Content models

- adapt content model of target XSD to information present in corpus = baseline
- compare derived content model with baseline
- $C_{XSD}, k = 2$ :
  - 38/47 as good
  - 9/47 better than baseline

ToSORE generalization + REDUCE smoothing

## 2/3. False positives/negatives

- $C_{XSD}$ ,  $k = 2$ 
  - $iXSD$ : no false positives/negatives
  - $iLOCAL$ , no REDUCE: 29 false positives

$$iXSD(k, C) = \text{REDUCE}(iLOCAL(k, C))$$

illustrates need for and power of REDUCE

- $C_1, \dots, C_8$ : no false positives/negatives

# Sensitivity to parameters $k$ and $\varepsilon$

- context size  $k \nearrow \Rightarrow$  false positives  $\nearrow$   
 $\Rightarrow$  false negatives  $\searrow$

rule of thumb:  
increase  $k$  until types are derived  
with too few examples

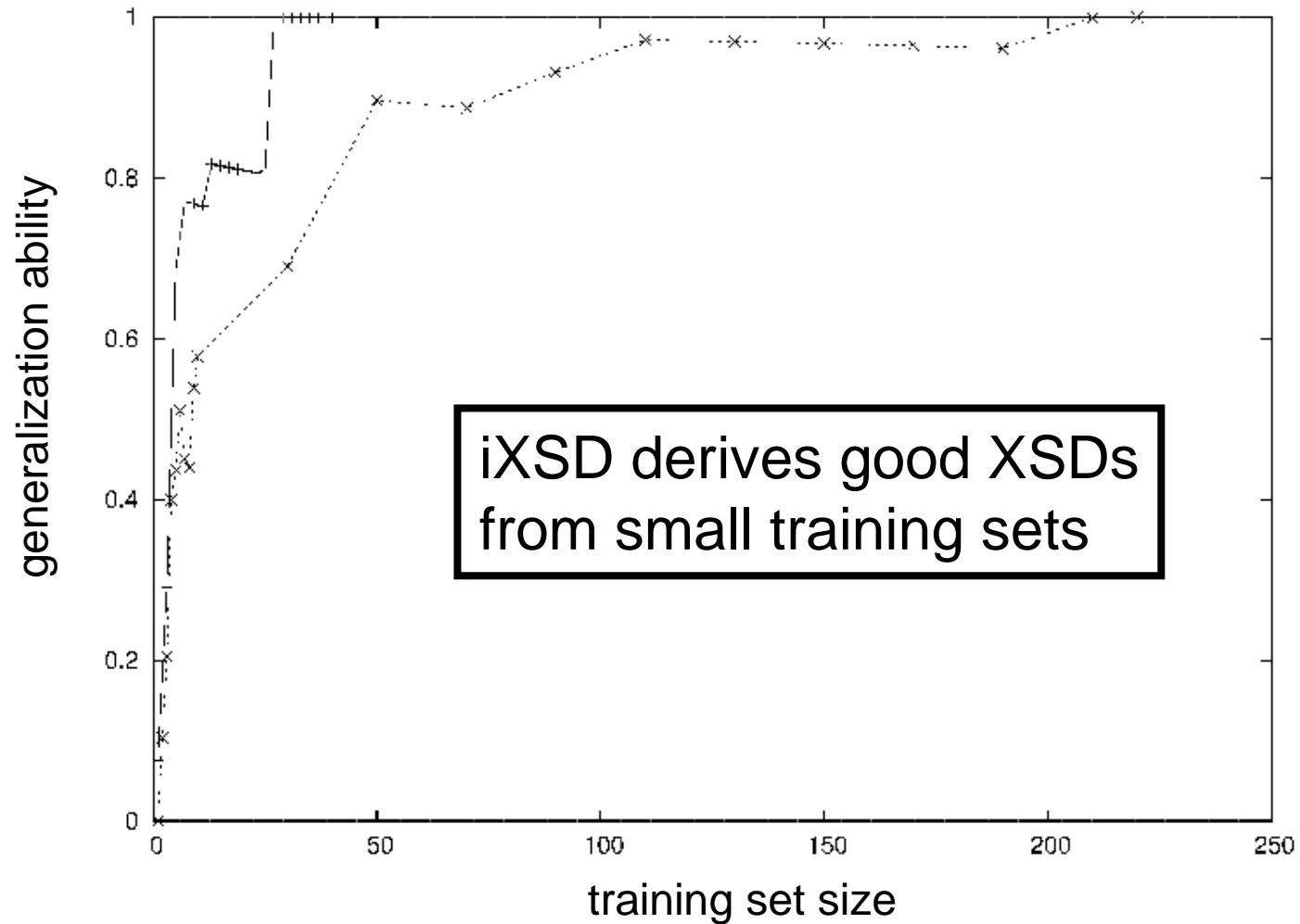
- $\varepsilon \nearrow \Rightarrow$  false positives  $\searrow$   
 $\Rightarrow$  false negatives  $\nearrow$

safe range:  $0.05 \lesssim \varepsilon \lesssim 0.15$

# Generalization

*i*XSD on training set

generalization ability = fraction of valid XML docs in test set



# Overview

- Introduction
- Complete algorithm *i*LOCAL
- Heuristic *i*XSD
- Experiments
- **Conclusions**

# Conclusions

- Two algorithms
  - *i*LOCAL: sound &  $k$ -complete
  - *i*XSD: extends *i*LOCAL to deal with poor data
    - good performance on real world & synthetic data
    - good runtime performance
    - rule of thumb to determine context size  $k$
- Future work
  - principled approach to determine best locality  $k$

Thank you