

# EXTRACTING DATA FROM THE WEB

Georg Gottlob

Oxford University

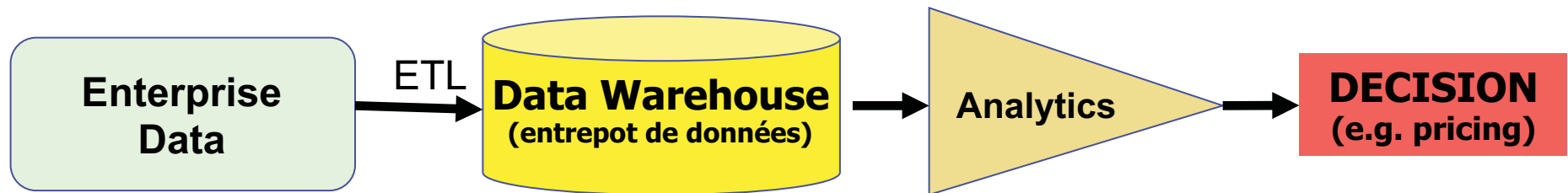


# Talk Outline

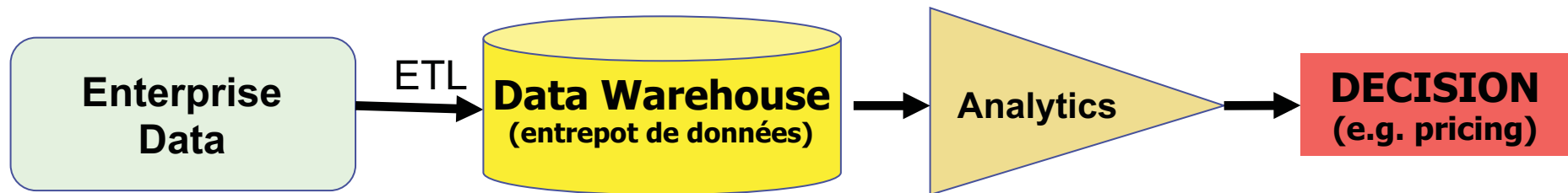
---

- Motivation: need of information extraction
- Logical foundations of information extraction
- The Lixto Visual Wrapper
- The Diadem Project

## Traditional data-based decision making in enterprises.



## Traditional data-based decision making in enterprises.



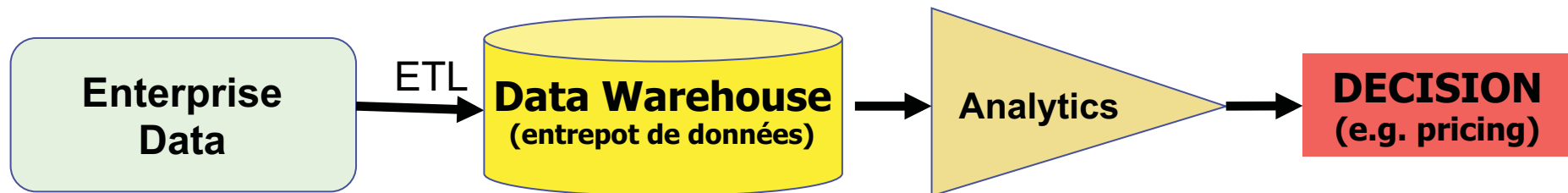
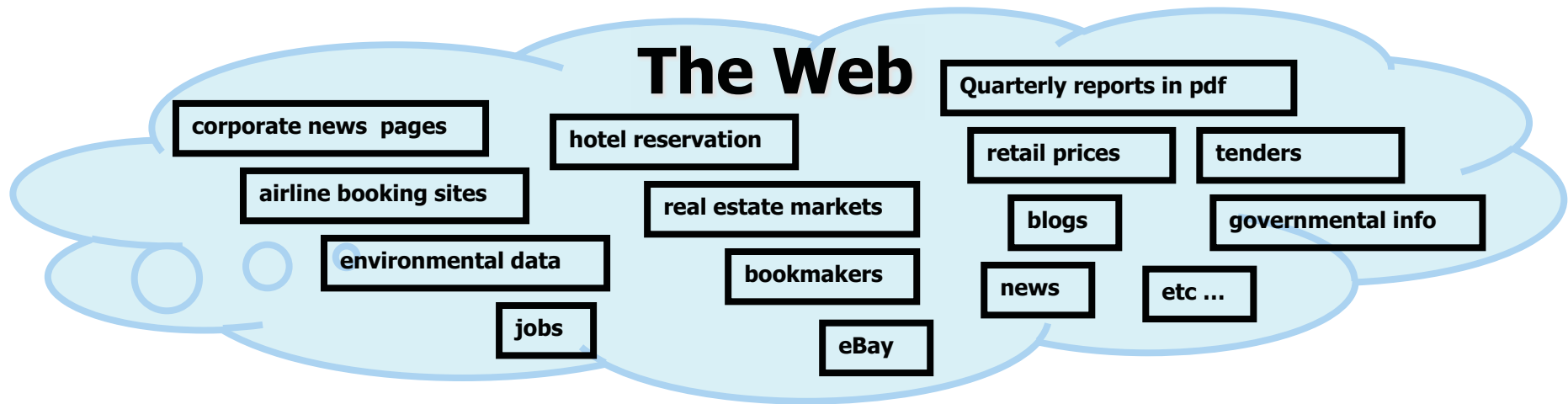
But often the most relevant data are outside the company, **on the Web!**

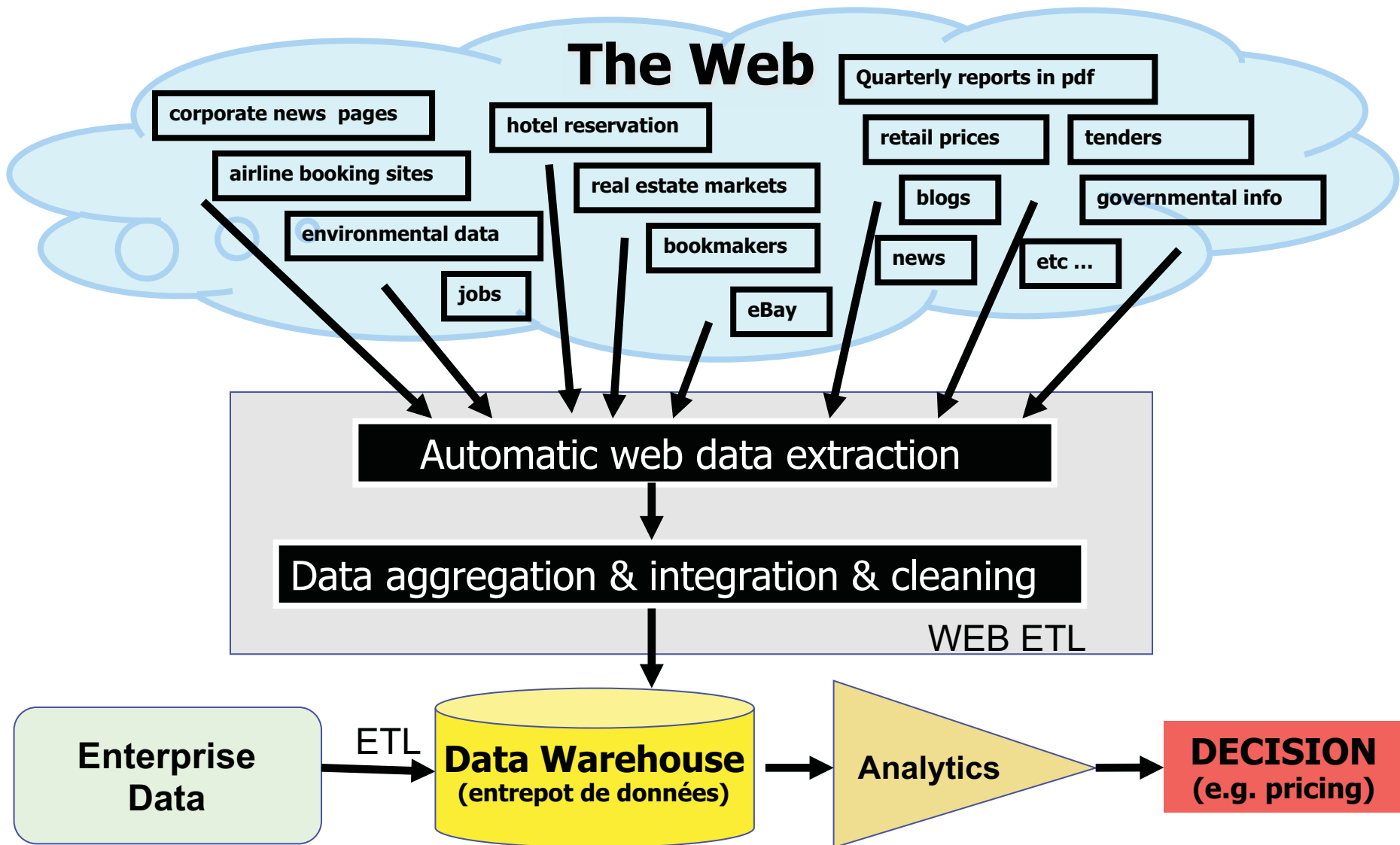
→ Online data intelligence, online market intelligence, automatic web data extraction.

# Online Market Intelligence (OMI)

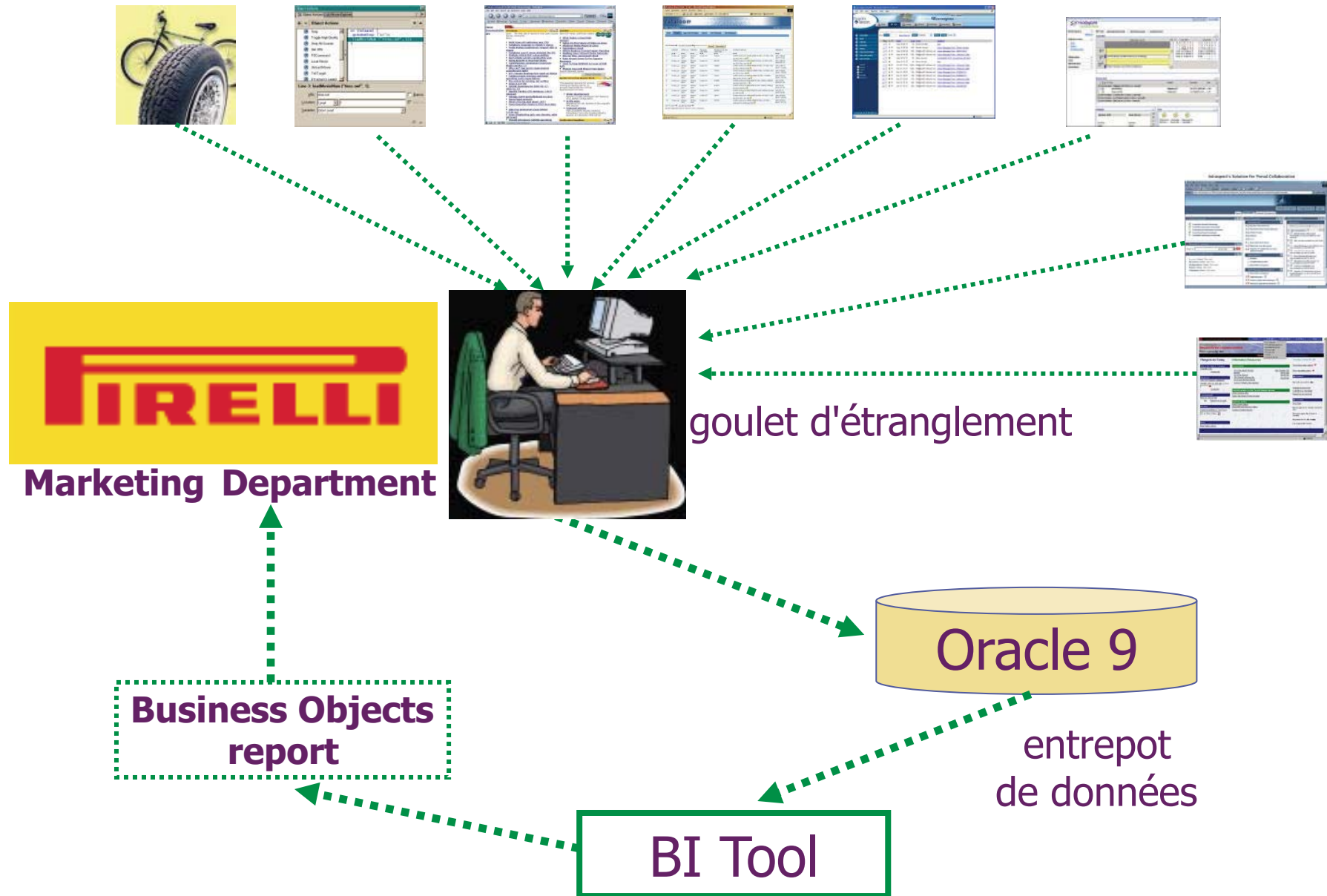
(surveillance du marché)

- **Electronics Retailer** (détaillant d'électronique - composants) :  
market overview, 20 competitors, 200,000 products/prices
- **Supermarket Chain:** Price comparison; must quickly react to  
special offers (offres spéciales) , new products,...
- **Internet Travel Agency:** Gives best price guarantee, wants to detect  
“pricing attacks”,...
- **Road Construction Company:** Find new public tenders (“appels d'offre”)
- **Hedge Fund** (“fonds de placement” ): Obtain recent house  
price changes from real-estate agent's  
Web pages before the weekly index is published. Anticipating the  
Consumer price index (index des prix à la consommation).
- **Governmental/Policy Making ....**





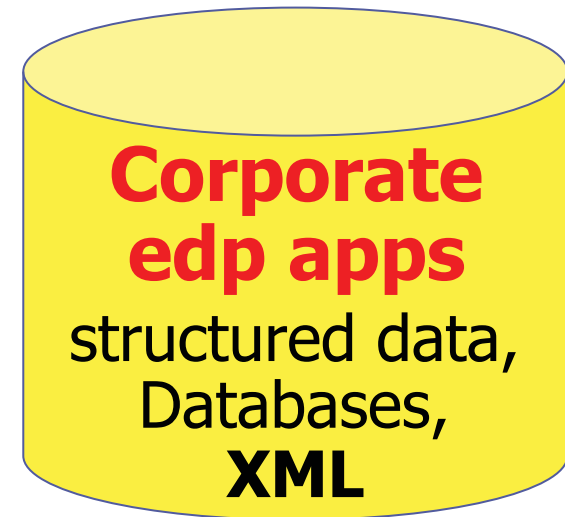
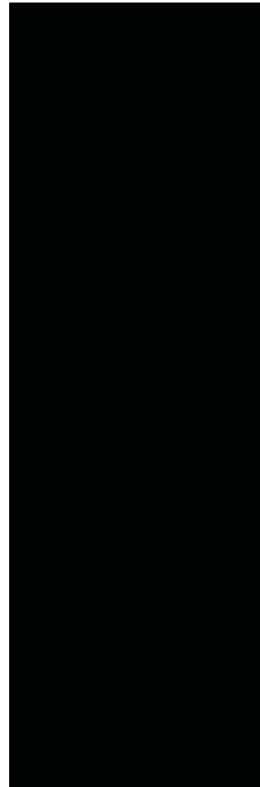
# Marketing & Business Intelligence

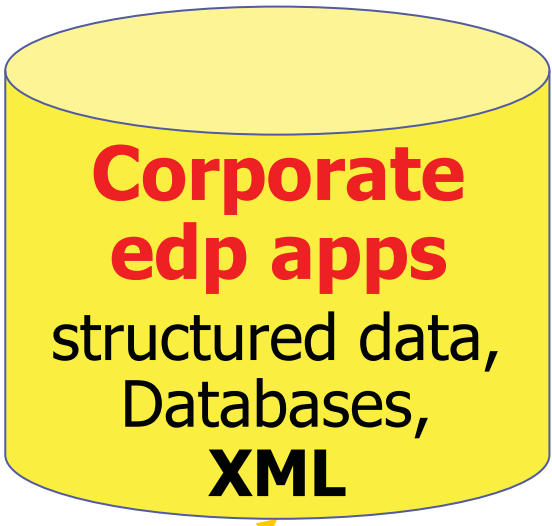
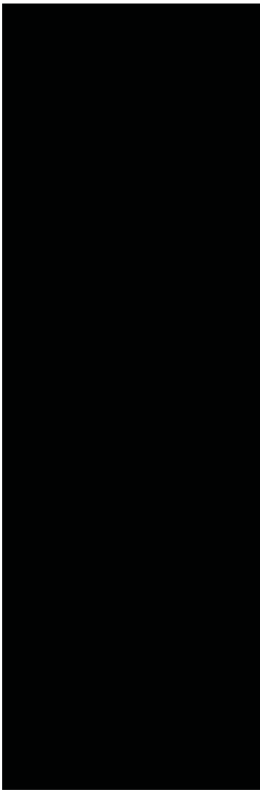




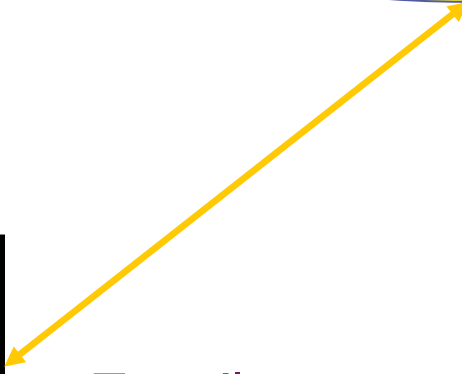
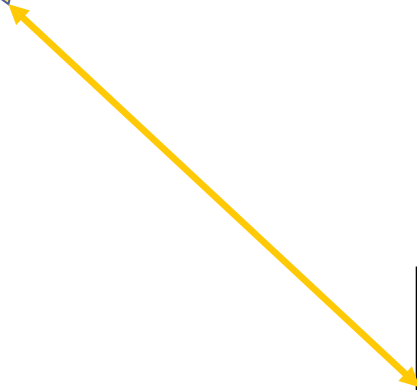
# The Wall

**Problem: Make web contents accessible to electronic data processing**



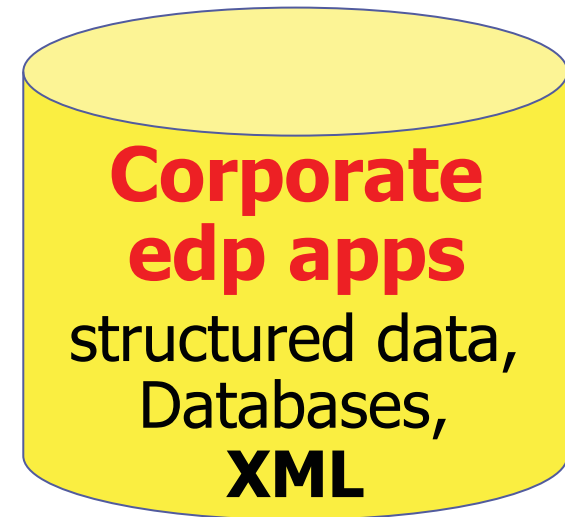


Travail  
aliénant



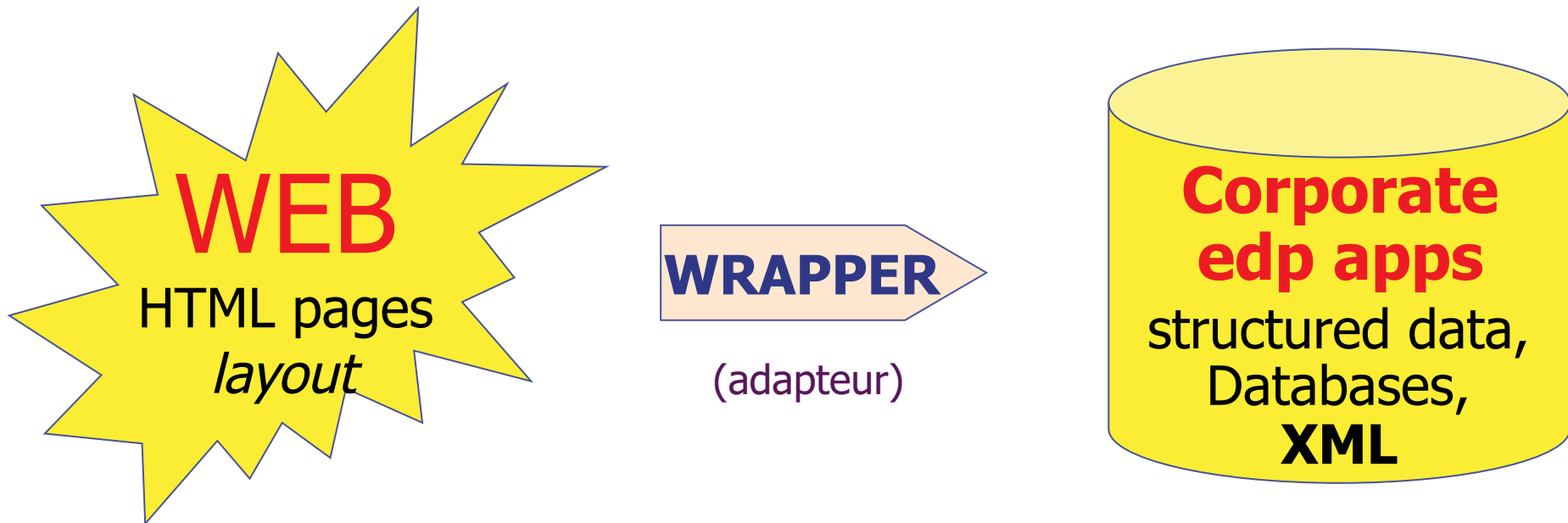
# Web wrapping

**Goal:** Make web contents accessible to electronic data processing



# Web wrapping

**Goal:** Make web contents accessible to electronic data processing



Wrappers: **HTML** → select → extract → annotate → **XML**

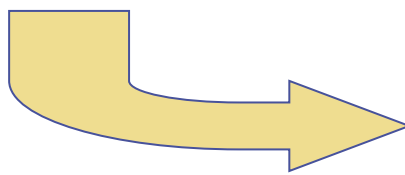
Home	Our offer of <b>195/65 R15 H</b>		
Advice	<a href="#">brand</a>	<a href="#">profile</a>	<a href="#">size</a> <a href="#">speed</a> <a href="#">price</a>
FAQs	Maximum 10 tyres will be displayed		
Customer comments	<b>Goodyear***</b>	<a href="#">EAGLE NCT 5</a>	With <b>mytyres.co.uk</b> only £ <b>44,10</b>
Help	Summer tyres	<a href="#">195/65 R15 91V</a>	<a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
Delivery and Payment	The proposed tyres could achieve higher speeds than those you are searching for. You can use these tyres without any hesitation		
Terms and Conditions	<b>Dunlop***</b>	<a href="#">SP SPORT 01</a>	With <b>mytyres.co.uk</b> only £ <b>44,30</b>
About us	Summer tyres	<a href="#">195/65 R15 91H</a>	<a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
My orders	<b>Pirelli***</b>	<a href="#">P 6000 Powergy</a>	With <b>mytyres.co.uk</b> only £ <b>44,70</b>
Newsletter	Summer tyres	<a href="#">195/65 R15 91H</a>	<a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
Tyre-Test	<b>Continental***</b>	<a href="#">EcoContact CP</a>	With <b>mytyres.co.uk</b> only £ <b>46,00</b>
Links	Summer tyres	<a href="#">195/65 R15 91H runout,</a>	<a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
Contact us	<b>Bridgestone***</b>	<a href="#">Turanza ER 31</a>	With <b>mytyres.co.uk</b> only £ <b>46,50</b>
	Summer tyres	<a href="#">195/65 R15 91H</a>	<a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
brought to you by <b>DELTA COM</b>	<b>Michelin***</b>	<a href="#">Pilot PRIMACY</a>	With <b>mytyres.co.uk</b> only £ <b>56,40</b>
	Summer tyres	<a href="#">195/65 R15 91H</a>	<a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
	<b>Goodyear***</b>	<a href="#">EAGLE NCT 5</a>	With <b>mytyres.co.uk</b> only £ <b>44,20</b>
	Summer tyres	<a href="#">195/65 R15 91H</a>	<a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
	<b>Pirelli***</b>	<a href="#">P 6000 Powergy</a>	With <b>mytyres.co.uk</b> only £ <b>44,70</b>
	Summer tyres	<a href="#">195/65 R15 91H</a>	<a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
	<b>Dunlop***</b>	<a href="#">SP SPORT 200</a>	With <b>mytyres.co.uk</b> only £ <b>46,40</b>
	Summer tyres	<a href="#">195/65 R15 91H E</a>	<a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
	<b>Pirelli***</b>	<a href="#">P 6</a>	With <b>mytyres.co.uk</b> only £ <b>44,70</b>
	Summer tyres	<a href="#">195/65 R15 91H</a>	<a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>

10 tyres from 118 were displayed. (1 - 10)

Page: [\[1\]](#) [\[2\]](#) [\[3\]](#) [\[4\]](#) [\[5\]](#) [\[6\]](#) [\[7\]](#) [\[8\]](#) [\[9\]](#) [\[10\]](#) [\[11\]](#) [\[12\]](#)

Prices includes postage, packing and VAT within mainland UK.

\*\*\*Please note: these tyres are subject to a delivery period of up to 2 weeks





Enregistrement:  
hierarchie de données

```
<?xml version="1.0" encoding="UTF-8"?>
<document>
  <tyre>
    <brand>Goodyear</brand>
    <profile>EAGLE NCT 5</profile>
    <price>44,10</price>
  </tyre>
  <tyre>
    <brand>Dunlop</brand>
    <profile>SP SPORT 01</profile>
    <price>44,30</price>
  </tyre>

```

# Patterns:

**mytyres.co.uk** The home of low tyre prices for cars, 4x4's and vans  
...fitting throughout the UK

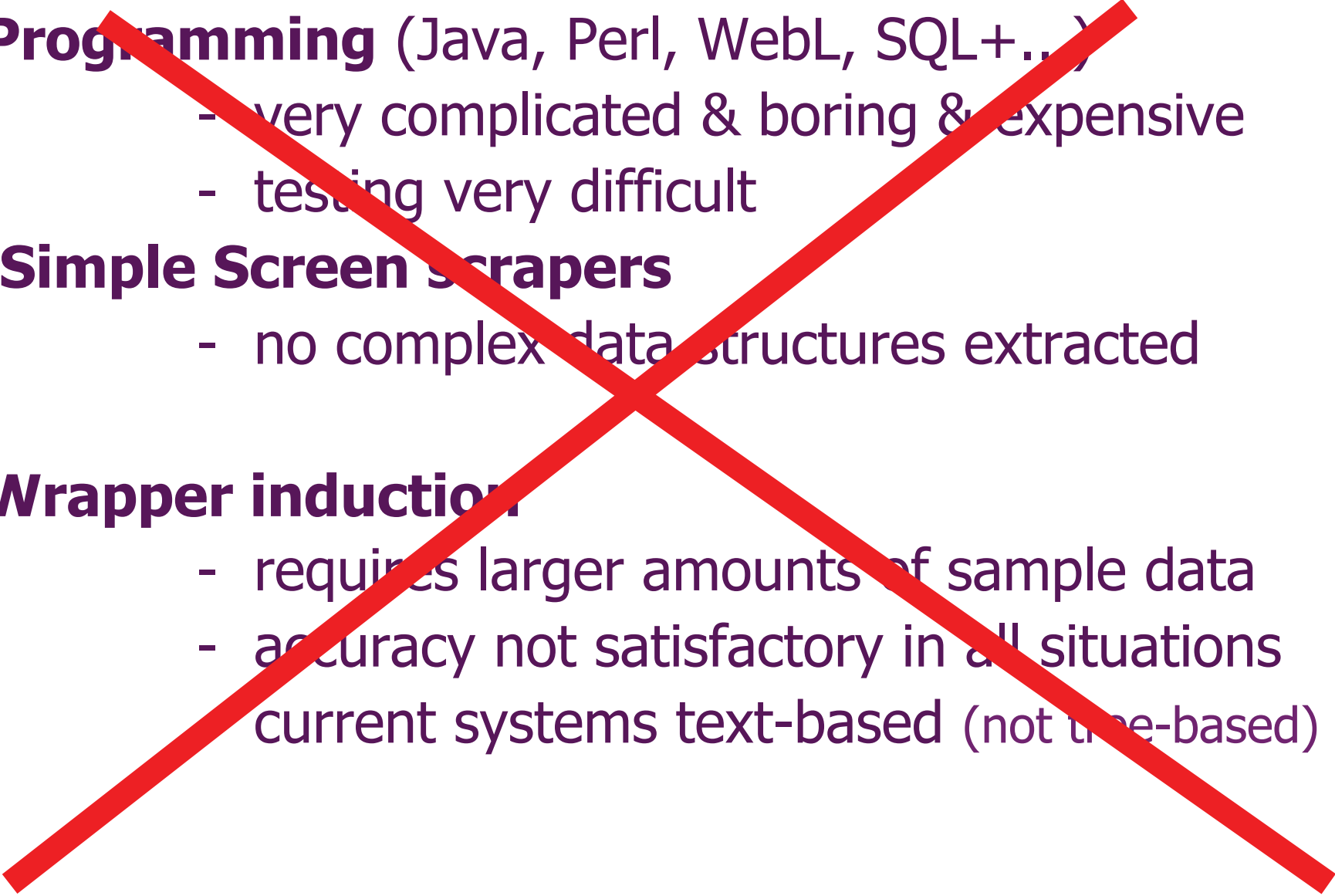
Home	Our offer of <i>195/65 R15 H</i>		
Advice	<a href="#">brand</a>	<a href="#">profile</a> <a href="#">size</a> <a href="#">speed</a>	<a href="#">price</a>
FAQs	Maximum 10 tyres will be displayed		
Customer comments	<b>Goodyear***</b>	<a href="#">EAGLE NCT 5</a> 195/65 R15 91V	With <b>mytyres.co.uk</b> only £ <b>44,10</b> <a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
Help	The proposed tyres could achieve higher speeds than those you are searching for. You can use these tyres without any hesitation		
Delivery and Payment	<b>Dunlop***</b>	<a href="#">SP SPORT 01</a> 195/65 R15 91H	With <b>mytyres.co.uk</b> only £ <b>44,30</b> <a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
Terms and Conditions	<b>Pirelli***</b>	<a href="#">P 6000 Powergy</a> 195/65 R15 91H	With <b>mytyres.co.uk</b> only £ <b>44,70</b> <a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
About us	<b>Continental***</b>	<a href="#">EcoContact CP</a> 195/65 R15 91H runout,	With <b>mytyres.co.uk</b> only £ <b>46,00</b> <a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
My orders	<b>Bridgestone***</b>	<a href="#">Turanza ER 31</a> 195/65 R15 91H	With <b>mytyres.co.uk</b> only £ <b>46,50</b> <a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
Newsletter	<b>Michelin***</b>	<a href="#">Pilot PRIMACY</a> 195/65 R15 91H	With <b>mytyres.co.uk</b> only £ <b>56,40</b> <a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
Tyre-Test	<b>Goodyear***</b>	<a href="#">EAGLE NCT 5</a> 195/65 R15 91H	With <b>mytyres.co.uk</b> only £ <b>44,20</b> <a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
Links	<b>Pirelli***</b>	<a href="#">P 6000 Powergy</a> 195/65 R15 91H	With <b>mytyres.co.uk</b> only £ <b>44,70</b> <a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
Contact us	<b>Dunlop***</b>	<a href="#">SP SPORT 200</a> 195/65 R15 91H E	With <b>mytyres.co.uk</b> only £ <b>46,40</b> <a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
	<b>Pirelli***</b>	<a href="#">P 6</a> 195/65 R15 91H	With <b>mytyres.co.uk</b> only £ <b>46,70</b> <a href="#">Details</a> <a href="#">Save for later</a> <a href="#">Buy now</a>
brought to you by 	10 tyres from 118 were displayed. (1 - 10)		
Page: <a href="#">[1]</a> <a href="#">[2]</a> <a href="#">[3]</a> <a href="#">[4]</a> <a href="#">[5]</a> <a href="#">[6]</a> <a href="#">[7]</a> <a href="#">[8]</a> <a href="#">[9]</a> <a href="#">[10]</a> <a href="#">[11]</a> <a href="#">[12]</a> <a href="#">next&gt;&gt;</a>			
Prices includes postage, packing and VAT within mainland UK.			
***Please note: these tyres are subject to a delivery period of up to 1 week.			

**Tyre** → **Brand** → **Profile** → **Price**

# Different approaches in the past

- ◆ **Programming** (Java, Perl, WebL, SQL+...)
  - very complicated & boring & expensive
  - testing very difficult
- ◆ **Simple Screen scrapers** ("gratte-écran")
  - no complex data structures extracted
- ◆ **Wrapper induction** (apprentissage d'adapteurs)
  - requires larger amounts of sample data
  - precision often not satisfactory
  - current systems text-based (not tree-based)

# Different approaches in the past

- ◆ **Programming** (Java, Perl, WebL, SQL+...)
    - very complicated & boring & expensive
    - testing very difficult
  - ◆ **Simple Screen scrapers**
    - no complex data structures extracted
  - ◆ **Wrapper induction**
    - requires larger amounts of sample data
    - accuracy not satisfactory in all situations
    - current systems text-based (not tree-based)
- 



# Modern Solutions

## ◆ **Semi-automatic tool (outils)**



- based on solid theory
- modular knowledge representation
- easy to use
- commercial product since 2002

## ◆ **Fully automated extraction**



- for specific application domains
- extracts from 1000s of websites
- current research

# Talk Outline

---

- Motivation: need of information extraction
- Logical foundations of information extraction
- The Lixto Visual Wrapper
- The Diadem Project

# Web documents are trees !

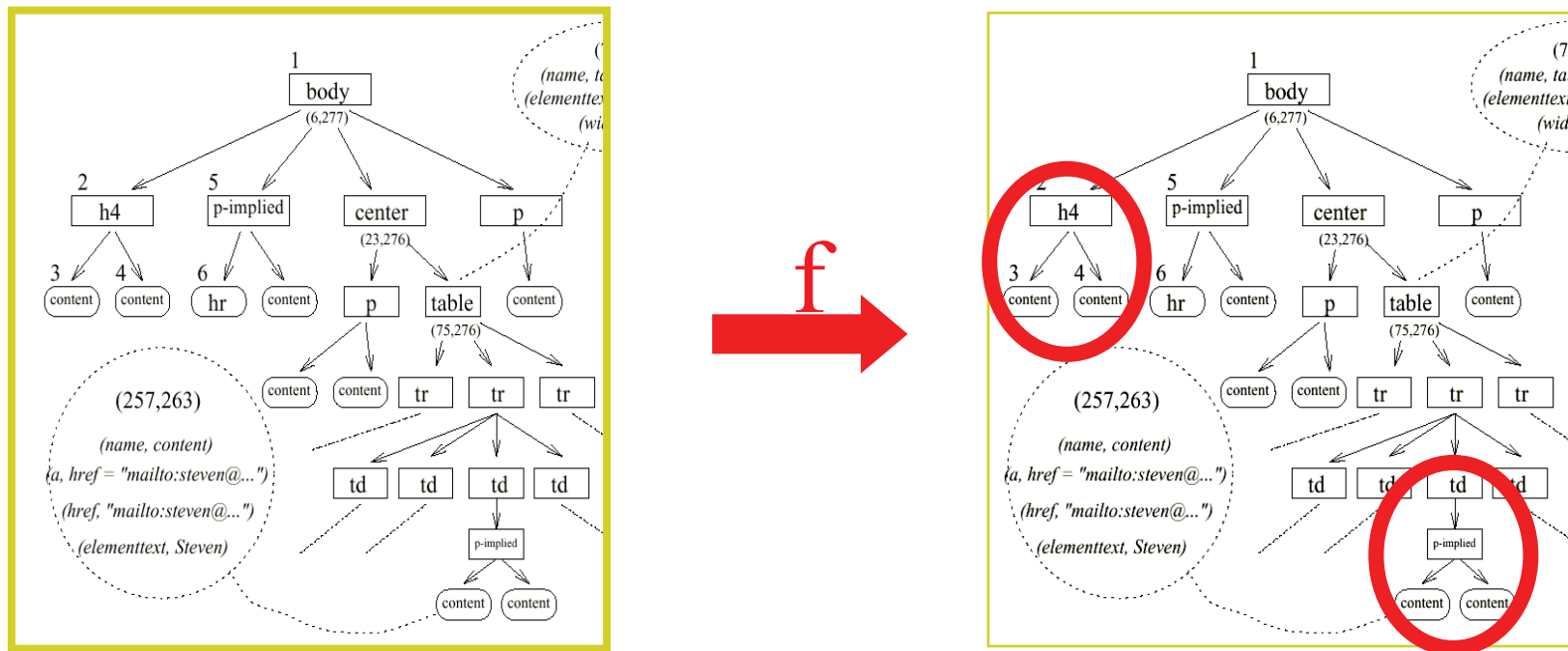
**HTML:** Hypertext Markup Language

**XML:** Extensible Markup Language

HTML, XML: Context free\* languages. Represent a document by its parse tree (arbre syntaxique).

# HTML Content Extractor

Function  $f$ : HTML Parse tree  $\rightarrow$  Subtrees



Leaves of subtrees are among leaves of orig. tree

# The Essence of Web Wrapping ?

Functional view: Wrapper defines functions  $f$

$$f: \text{Tree} \rightarrow \mathcal{P}(\text{Tree})$$

$$t \rightarrow T \subseteq \text{subtrees}(t)$$

Equivalent logical view:

**Wrapper defines monadic predicates  $P$   
over the nodes (arbre dom) of each input  
document**

# A HTML page

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN">
```

```
<html> <body>
```

```
<h1>People @ DBAI</h1>
```

```
<table border="1" cellpadding="3" cellspacing="1">
```

```
  <tr> <td>Georg Gottlob</td>
```

```
    <td>gottlob@dbai.tuwien.ac.at</td>
```

```
    <td>18420</td>
```

```
</tr>
```

```
<tr> <td>Christoph Koch</td>
```

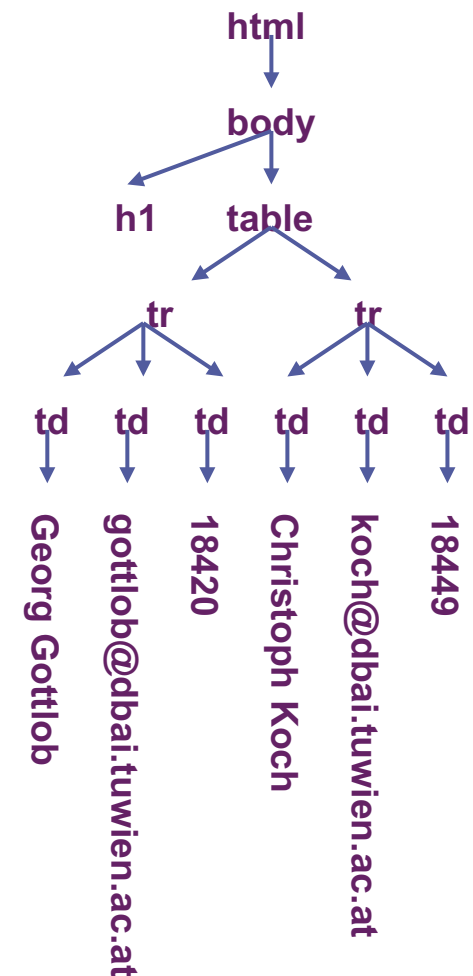
```
  <td>koch@dbai.tuwien.ac.at</td>
```

```
  <td>18449</td>
```

```
</tr>
```

```
</table>
```

```
</body> </html>
```



## People @ DBAI

Georg Gottlob	gottlob@...	18420
Christoph Koch	koch@...	18449

# Predicate *employeetable*

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN">
```

```
<html> <body>
```

```
<h1>People @ DBAI</h1>
```

```
<table border="1" cellpadding="3" cellspacing="1">
```

```
<tr> <td>Georg Gottlob</td>
```

```
<td>gottlob@dbai.tuwien.ac.at</td>
```

```
<td>18420</td>
```

```
</tr>
```

```
<tr> <td>Christoph Koch</td>
```

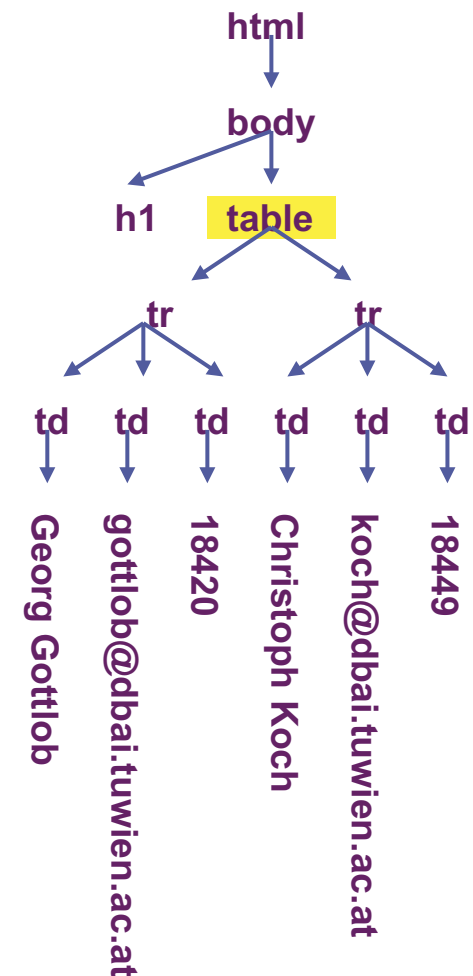
```
<td>koch@dbai.tuwien.ac.at</td>
```

```
<td>18449</td>
```

```
</tr>
```

```
</table>
```

```
</body> </html>
```



## People @ DBAI

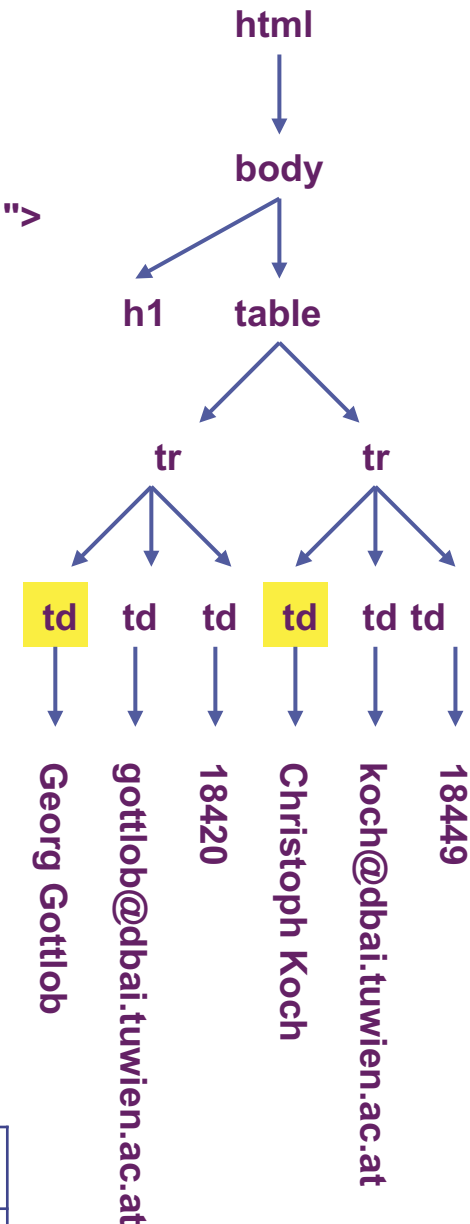
Georg Gottlob	gottlob@...	18420
Christoph Koch	koch@...	18449

# Predicate *employee*

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN">
<html> <body>
<h1>People @ DBAI</h1>
<table border="1" cellpadding="3" cellspacing="1">
  <tr> <td>Georg Gottlob</td>
    <td>gottlob@dbai.tuwien.ac.at</td>
    <td>18420</td>
  </tr>
  <tr> <td>Christoph Koch</td>
    <td>koch@dbai.tuwien.ac.at</td>
    <td>18449</td>
  </tr>
</table>
</body> </html>
```

## People @ DBAI

Georg Gottlob	gottlob@...	18420
Christoph Koch	koch@...	18449



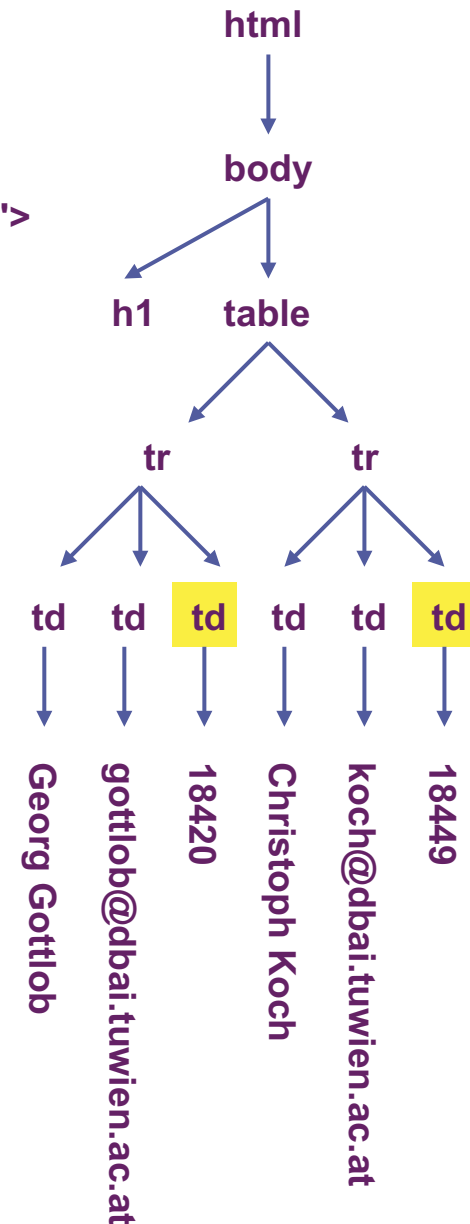


# Predicate *phone*

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN">
<html> <body>
<h1>People @ DBAI</h1>
<table border="1" cellpadding="3" cellspacing="1">
  <tr> <td>Georg Gottlob</td>
    <td>gottlob@dbai.tuwien.ac.at</td>
    <td>18420</td>
  </tr>
  <tr> <td>Christoph Koch</td>
    <td>koch@dbai.tuwien.ac.at</td>
    <td>18449</td>
  </tr>
</table>
</body> </html>
```

## People @ DBAI

Georg Gottlob	gottlob@...	18420
Christoph Koch	koch@...	18449



# Expressiveness Yardstick: MSO

- MSO captures exactly the essence of data extraction:
  - Define sets of nodes of a document
- Expressiveness, complexity, semantics well understood:
  - ✓ MSO over trees: perfect logical semantics
  - ✓ MSO over trees: high expressive power (tree automata)
  - ✓ MSO over trees: low data complexity
- Drawbacks:
  - hard to use, no visual specification,
  - high query complexity (cpl. de requetes)(→ bad scalability, mauvais passage à l' échelle).

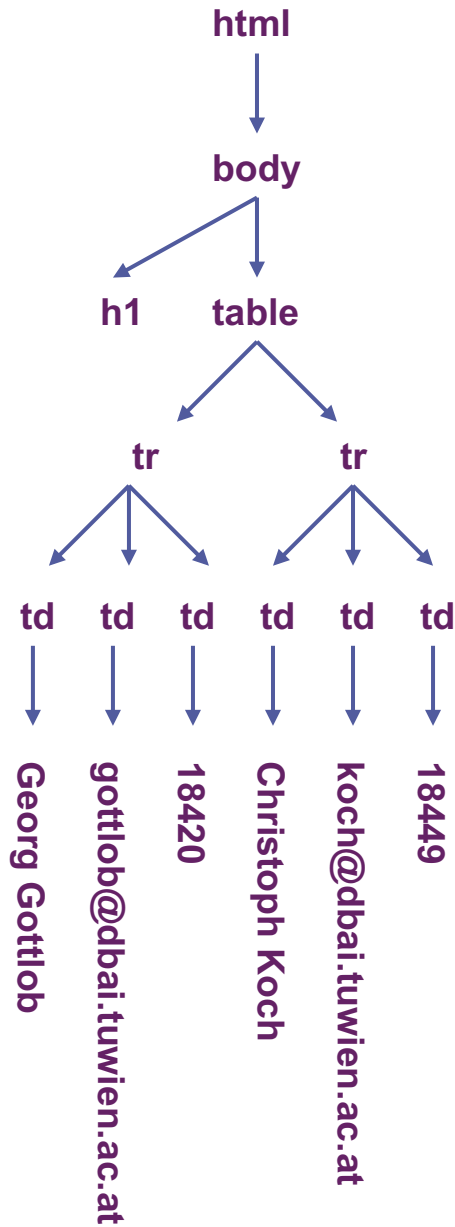
# MSO on strings and trees



## Rich theory:

- Büchi: MSO = REG over strings (chaînes de caractères)
- Thatcher and Wright, Rabin:  
MSO = REG over ranked trees (arbres bornés)  
= tree automata
- Brüggemann-Klein/Wood/Murata:  
MSO = REG over unranked trees
- Neven & Schwentick: Unranked Query Automata
- Courcelle: MSO in LinTime on tree-like structures  
(treewidth  $\leq k$ , data complexity)

# Ordered Trees as finite structures



**firstchild**

**nextsibling**

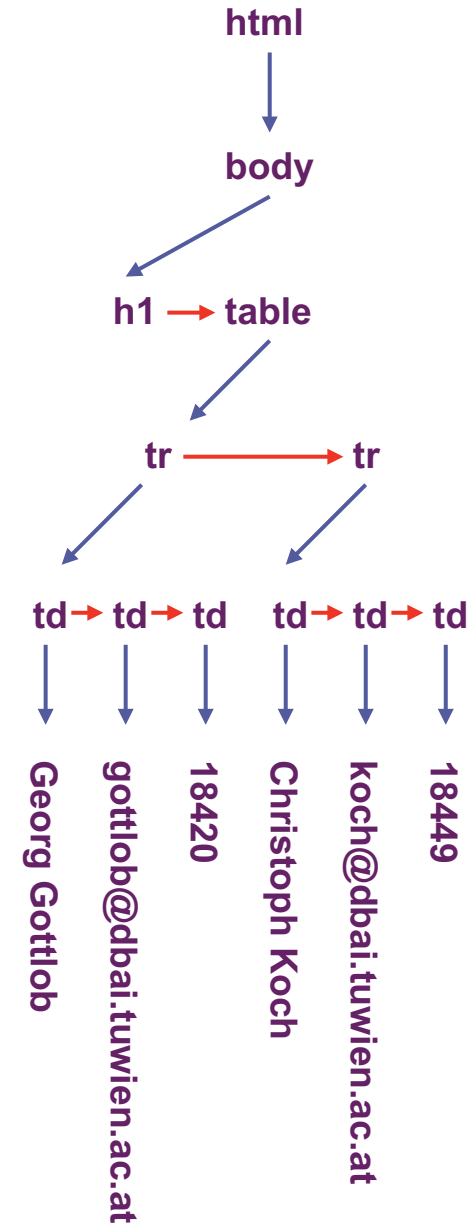
$\text{label}_{h1}()$

$\text{label}_{td}()$

...

$\text{root}()$

$\text{leaf}()$

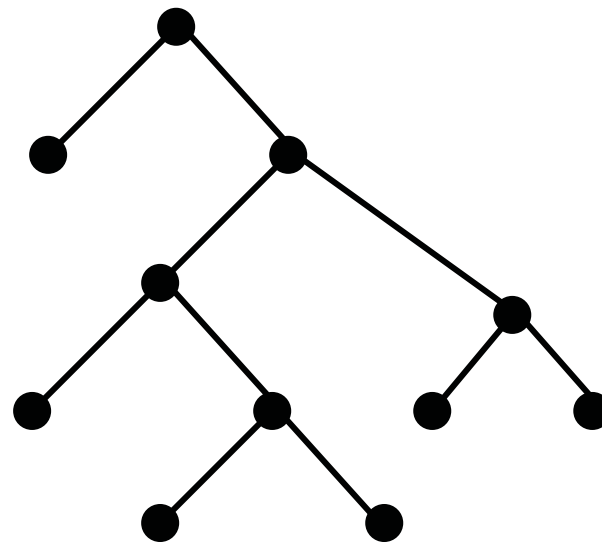


# MSO over Trees

Extract from a binary tree all roots of sub-trees with an odd number of leaves:

$$\exists S \forall x [ S(u) \ \& \ (\text{leaf}(x) \rightarrow S(x)) \ \& \\ \forall x, y, z (((\text{firstchild}(x, y) \ \& \ \text{nextsibling}(y, z)) \rightarrow \\ (S(x) \leftrightarrow \neg(S(y) \leftrightarrow S(z)))))]$$

Tree automaton:



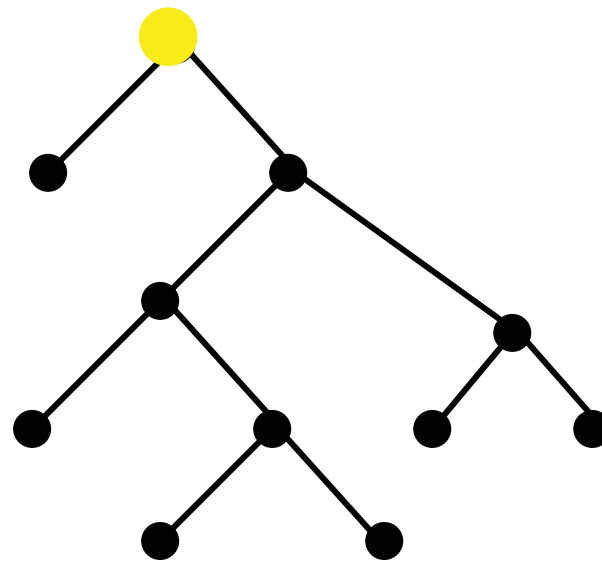
- auxiliary state
- roots even subtree
- roots odd subtree

# MSO over Trees

Extract from a binary tree all roots of sub-trees with an odd number of leaves:

$$\exists S \forall x [ S(x) \ \& \ (\text{leaf}(x) \rightarrow S(x)) \ \& \\ \forall x, y, z (((\text{firstchild}(x, y) \ \& \ \text{nextsibling}(y, z)) \rightarrow \\ (S(x) \leftrightarrow \neg(S(y) \leftrightarrow S(z)))))]$$

Tree automaton:

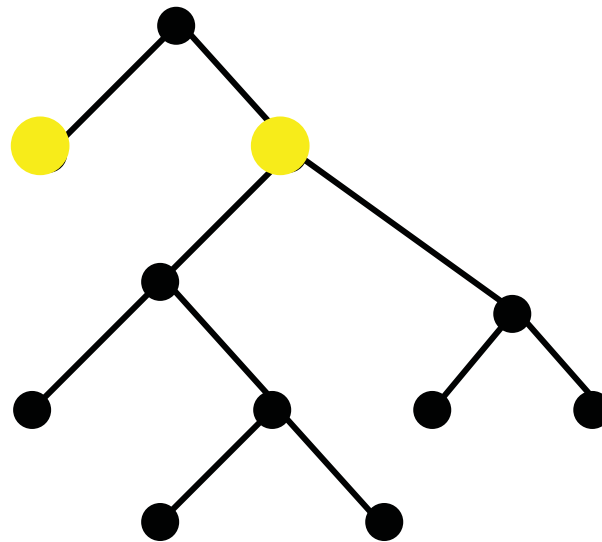


# MSO over Trees

Extract from a binary tree all roots of sub-trees with an odd number of leaves:

$$\exists S \forall x [ S(x) \ \& \ (\text{leaf}(x) \rightarrow S(x)) \ \& \\ \forall x, y, z (((\text{firstchild}(x, y) \ \& \ \text{nextsibling}(y, z)) \rightarrow \\ (S(x) \leftrightarrow \neg(S(y) \leftrightarrow S(z)))))]$$

Tree automaton:

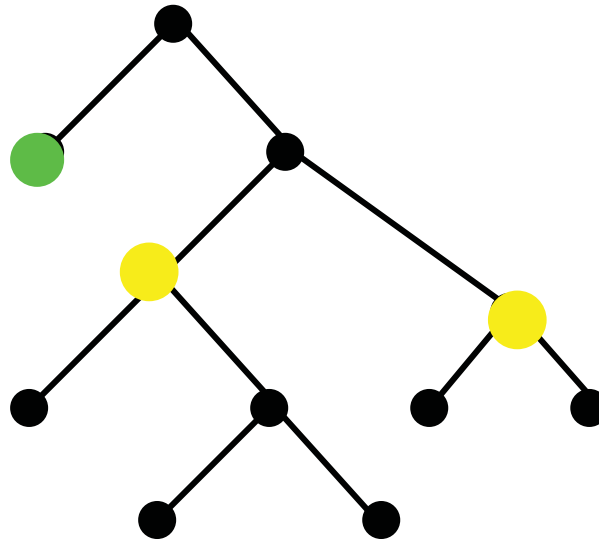


# MSO over Trees

Extract from a binary tree all roots of sub-trees with an odd number of leaves:

$$\exists S \forall x [ S(x) \ \& \ (\text{leaf}(x) \rightarrow S(x)) \ \& \\ \forall x, y, z (((\text{firstchild}(x, y) \ \& \ \text{nextsibling}(y, z)) \rightarrow \\ (S(x) \leftrightarrow \neg(S(y) \leftrightarrow S(z)))))]$$

Tree automaton:



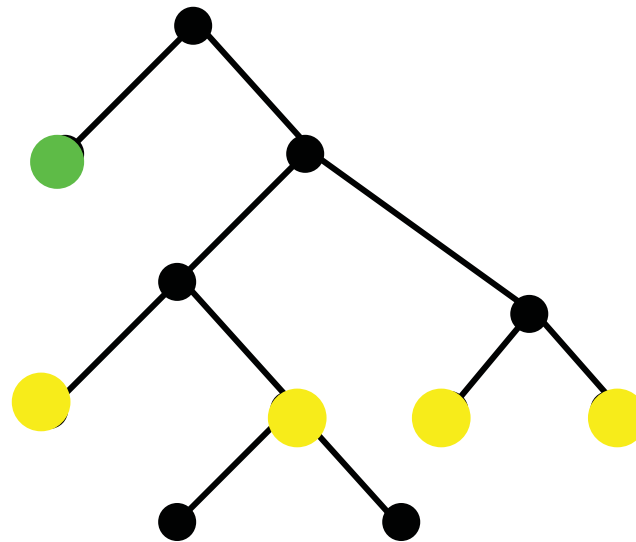


# MSO over Trees

Extract from a binary tree all roots of sub-trees with an odd number of leaves:

$$\exists S \forall x [ S(x) \ \& \ (\text{leaf}(x) \rightarrow S(x)) \ \& \\ \forall x, y, z ((\text{firstchild}(x, y) \ \& \ \text{nextsibling}(y, z)) \rightarrow \\ (S(x) \leftrightarrow \neg(S(y) \leftrightarrow S(z))))]$$

Tree automaton:

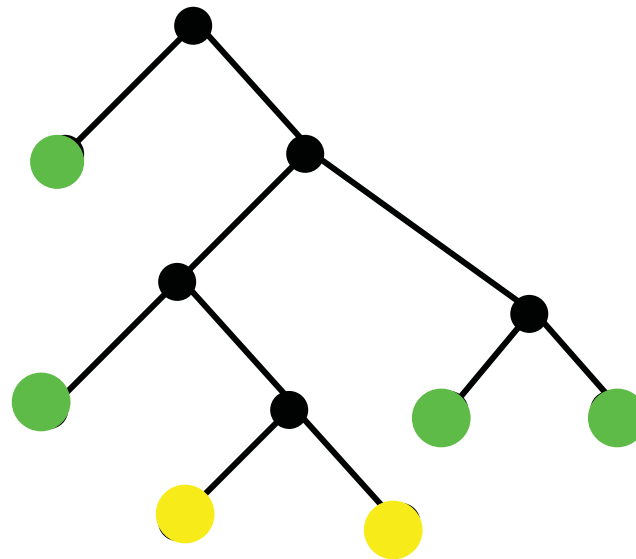


# MSO over Trees

Extract from a binary tree all roots of sub-trees with an odd number of leaves:

$$\exists S \forall x [ S(x) \ \& \ (\text{leaf}(x) \rightarrow S(x)) \ \& \\ \forall x, y, z (((\text{firstchild}(x, y) \ \& \ \text{nextsibling}(y, z)) \rightarrow \\ (S(x) \leftrightarrow \neg(S(y) \leftrightarrow S(z)))))]$$

Tree automaton:

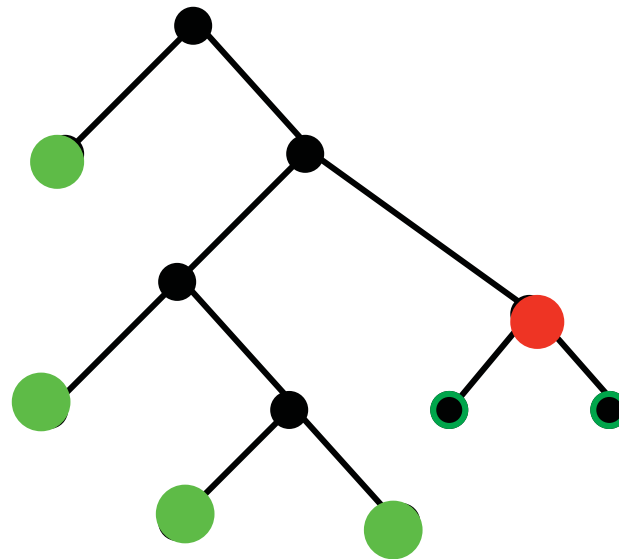


# MSO over Trees

Extract from a binary tree all roots of sub-trees with an odd number of leaves:

$$\exists S \forall x [ S(x) \ \& \ (\text{leaf}(x) \rightarrow S(x)) \ \& \\ \forall x, y, z (((\text{firstchild}(x, y) \ \& \ \text{nextsibling}(y, z)) \rightarrow \\ (S(x) \leftrightarrow \neg(S(y) \leftrightarrow S(z)))))]$$

Tree automaton:



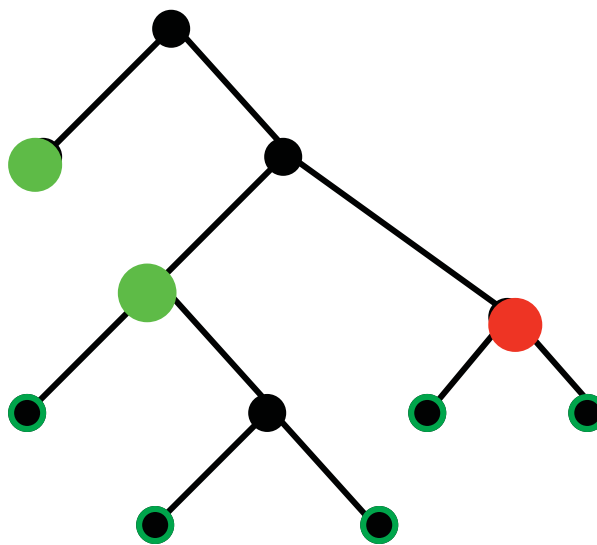


# MSO over Trees

Extract from a binary tree all roots of sub-trees with an odd number of leaves:

$$\begin{aligned} \exists S \forall x [ & S(x) \ \& \ (\text{leaf}(x) \rightarrow S(x)) \ \& \\ & \forall x, y, z (((\text{firstchild}(x, y) \ \& \ \text{nextsibling}(y, z)) \rightarrow \\ & (S(x) \leftrightarrow \neg(S(y) \leftrightarrow S(z)))))] \end{aligned}$$

Tree automaton:

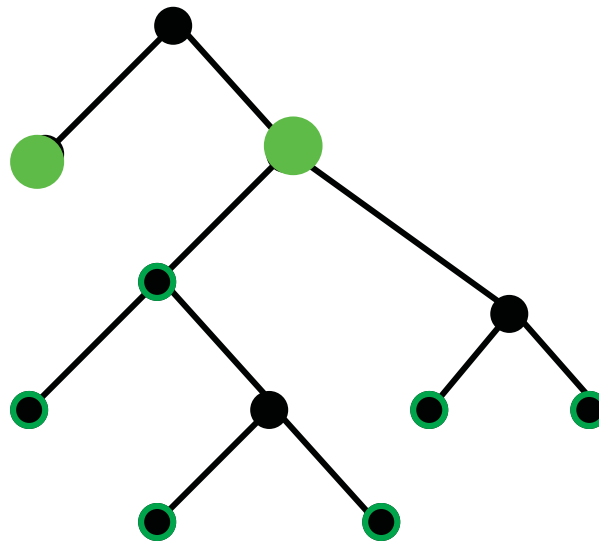


# MSO over Trees

Extract from a binary tree all roots of sub-trees with an odd number of leaves:

$$\exists S \forall x [ S(x) \ \& \ (\text{leaf}(x) \rightarrow S(x)) \ \& \\ \forall x, y, z (((\text{firstchild}(x, y) \ \& \ \text{nextsibling}(y, z)) \rightarrow \\ (S(x) \leftrightarrow \neg(S(y) \leftrightarrow S(z)))))]$$

Tree automaton:

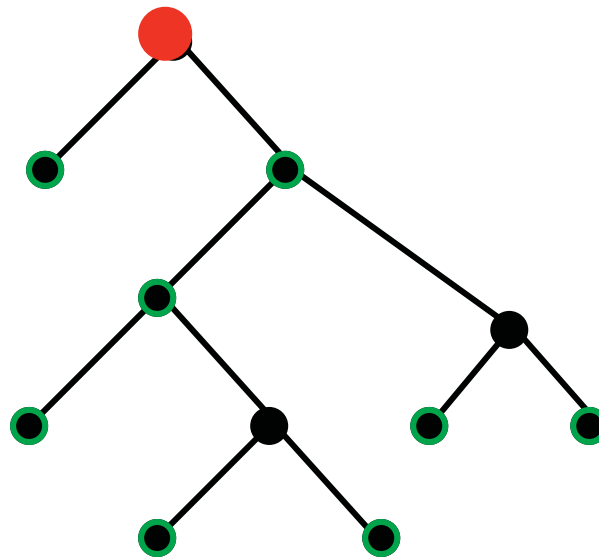


# MSO over Trees

Extract from a binary tree all roots of sub-trees with an odd number of leaves:

$$\exists S \forall x [ S(x) \ \& \ (\text{leaf}(x) \rightarrow S(x)) \ \& \\ \forall x, y, z (((\text{firstchild}(x, y) \ \& \ \text{nextsibling}(y, z)) \rightarrow \\ (S(x) \leftrightarrow \neg(S(y) \leftrightarrow S(z)))))]$$

Tree automaton:

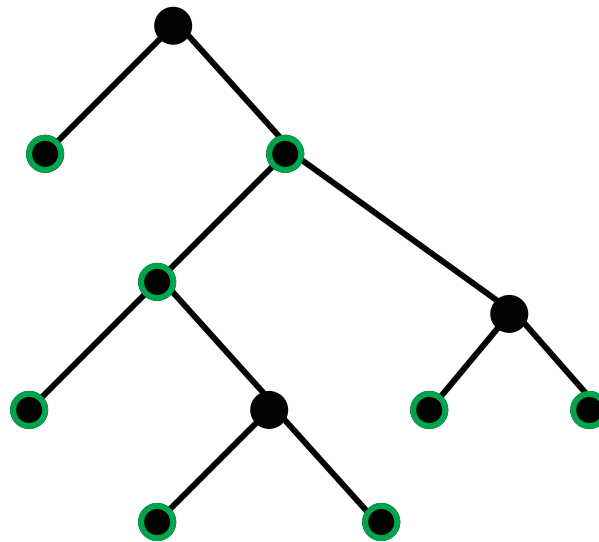


# MSO over Trees

Extract from a binary tree all roots of sub-trees with an odd number of leaves:

$$\exists S \forall x [ S(x) \ \& \ (\text{leaf}(x) \rightarrow S(x)) \ \& \\ \forall x, y, z (((\text{firstchild}(x, y) \ \& \ \text{nextsibling}(y, z)) \rightarrow \\ (S(x) \leftrightarrow \neg(S(y) \leftrightarrow S(z)))))]$$

Tree automaton:





**Logic heaven**

**DB theory heaven**

**DB programming  
heaven**

**Application design  
heaven**



**Logic heaven**

**DB theory heaven**

**DB programming  
heaven**

**Application design  
heaven**



**MSO**

**Logic heaven**

**DB theory heaven**

**DB programming  
heaven**

**Application design  
heaven**



**MSO**

**||**

**Monadic Datalog**

**Logic heaven**

**DB theory heaven**

**DB programming  
heaven**

**Application design  
heaven**



**MSO**

**II**

**Monadic Datalog**

**IN**

**Elog**

**Logic heaven**

**DB theory heaven**

**DB programming  
heaven**

**Application design  
heaven**



**MSO**

**||**

**Monadic Datalog**

**∩**

**Elog**

**∩**

**Lixto Visual Wrapper**

**Logic heaven**

**DB theory heaven**

**DB programming  
heaven**

**Application design  
heaven**



**MSO**

||

**Monadic Datalog**

∩

**Elog**

∩

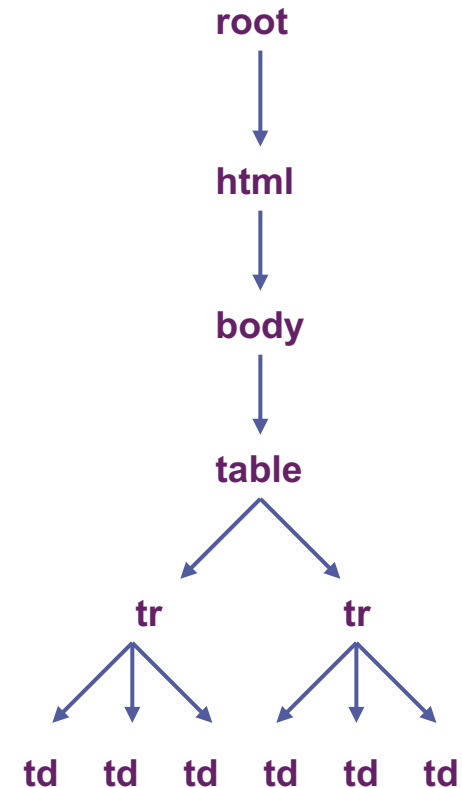
**Lixto Visual Wrapper**

∩

**liXto** Suite

# Monadic Datalog as a Wrapping Language

```
entry(X) :- root(R), firstchild(R,U), label[html](U),  
            firstchild(U,V), label[body](V),  
            firstchild(V,W),label[table](W),  
            firstchild(W,X), label[tr](X).  
entry(X):- entry(Y), nextsibling(Y,X).  
name(X) :- entry(E), firstchild(E, X), label[td](X).  
email(X) :- name(N), nextsibling(N, X), label[td](X).  
phone(X) :- email(M), nextsibling(M, X), label[td](X).
```



# Monadic Datalog as a Wrapping Language

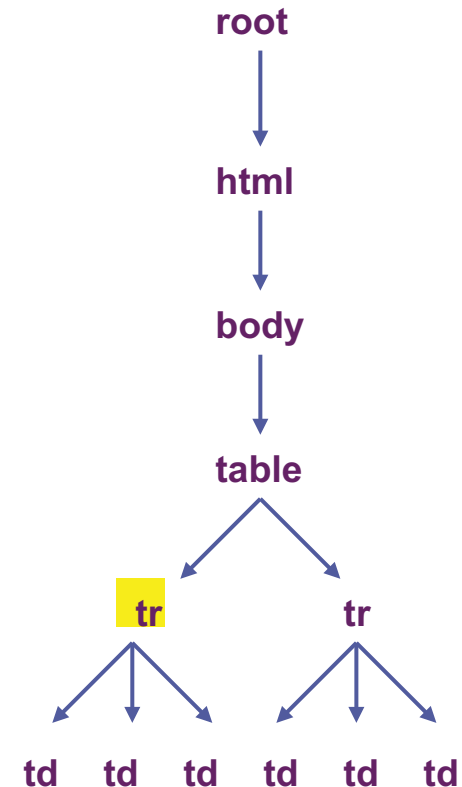
```
entry(X) :- root(R), firstchild(R,U), label[html](U),  
            firstchild(U,V), label[body](V),  
            firstchild(V,W), label[table](W),  
            firstchild(W,X), label[tr](X).
```

```
entry(X):- entry(Y), nextsibling(Y,X).
```

```
name(X) :- entry(E), firstchild(E, X), label[td](X).
```

```
email(X) :- name(N), nextsibling(N, X), label[td](X).
```

```
phone(X) :- email(M), nextsibling(M, X), label[td](X).
```





# Monadic Datalog as a Wrapping Language

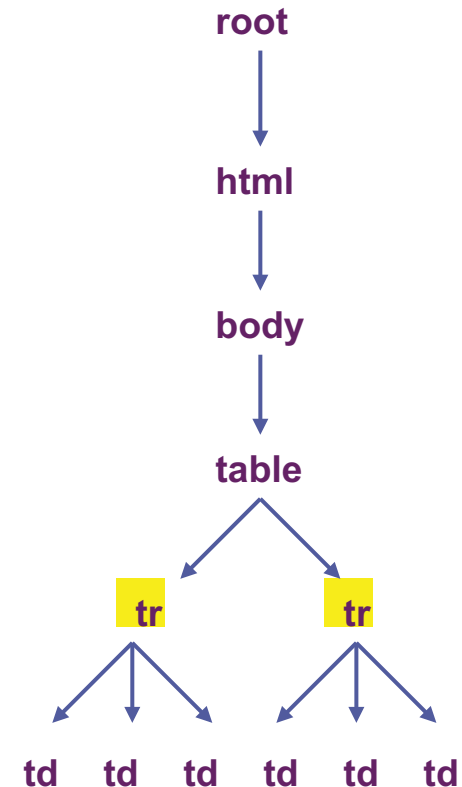
```
entry(X) :- root(R), firstchild(R,U), label[html](U),  
            firstchild(U,V), label[body](V),  
            firstchild(V,W), label[table](W),  
            firstchild(W,X), label[tr](X).
```

```
entry(X):- entry(Y), nextsibling(Y,X).
```

```
name(X) :- entry(E), firstchild(E, X), label[td](X).
```

```
email(X) :- name(N), nextsibling(N, X), label[td](X).
```

```
phone(X) :- email(M), nextsibling(M, X), label[td](X).
```



# Monadic Datalog as a Wrapping Language

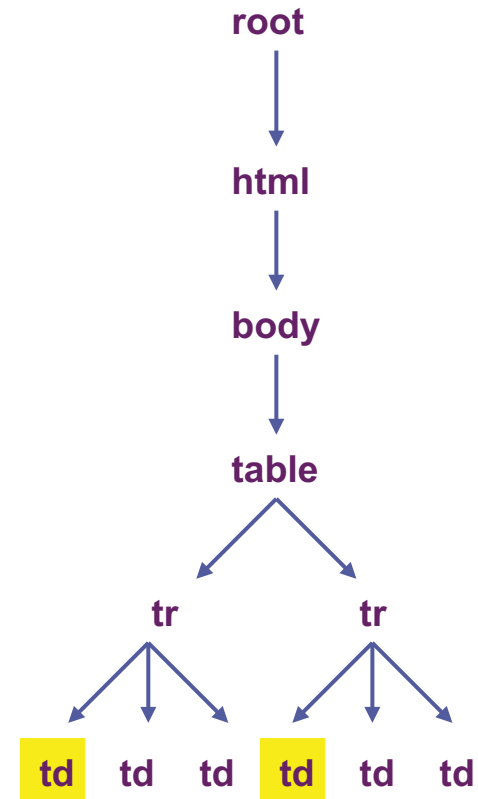
```
entry(X) :- root(R), firstchild(R,U), label[html](U),  
            firstchild(U,V), label[body](V),  
            firstchild(V,W), label[table](W),  
            firstchild(W,X), label[tr](X).
```

```
entry(X):- entry(Y), nextsibling(Y,X).
```

```
name(X) :- entry(E), firstchild(E, X), label[td](X).
```

```
email(X) :- name(N), nextsibling(N, X), label[td](X).
```

```
phone(X) :- email(M), nextsibling(M, X), label[td](X).
```



# Monadic Datalog as a Wrapping Language

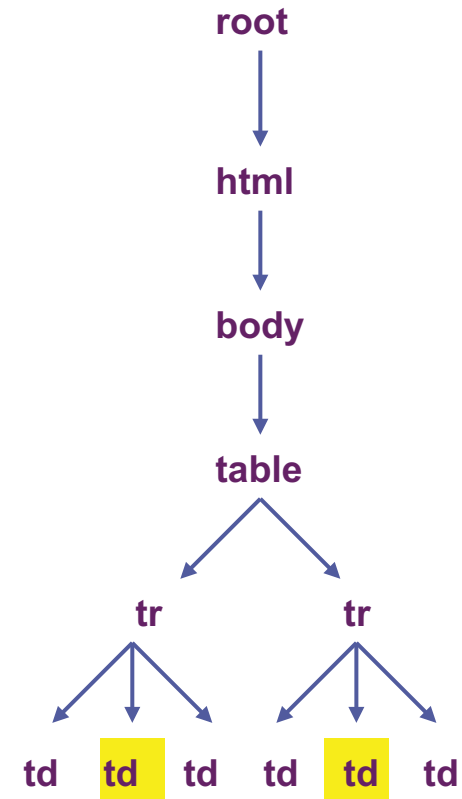
```
entry(X) :- root(R), firstchild(R,U), label[html](U),  
            firstchild(U,V), label[body](V),  
            firstchild(V,W), label[table](W),  
            firstchild(W,X), label[tr](X).
```

```
entry(X):- entry(Y), nextsibling(Y,X).
```

```
name(X) :- entry(E), firstchild(E, X), label[td](X).
```

```
email(X) :- name(N), nextsibling(N, X), label[td](X).
```

```
phone(X) :- email(M), nextsibling(M, X), label[td](X).
```



# Monadic Datalog as a Wrapping Language

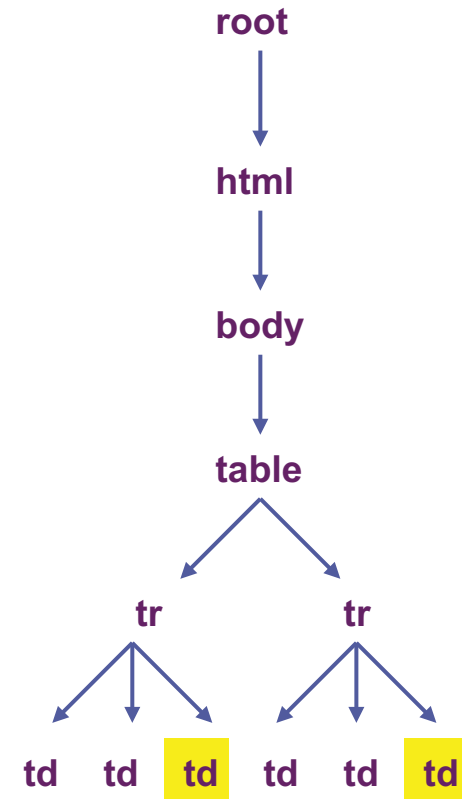
```
entry(X) :- root(R), firstchild(R,U), label[html](U),  
            firstchild(U,V), label[body](V),  
            firstchild(V,W), label[table](W),  
            firstchild(W,X), label[tr](X).
```

```
entry(X):- entry(Y), nextsibling(Y,X).
```

```
name(X) :- entry(E), firstchild(E, X), label[td](X).
```

```
email(X) :- name(N), nextsibling(N, X), label[td](X).
```

```
phone(X) :- email(M), nextsibling(M, X), label[td](X).
```



```
entry(X) :- root(R), firstchild(R,U), label[html](U),
            firstchild(U,V), label[body](V),
            firstchild(V,W),label[table](W),
            firstchild(W,X), label[tr](X).
```

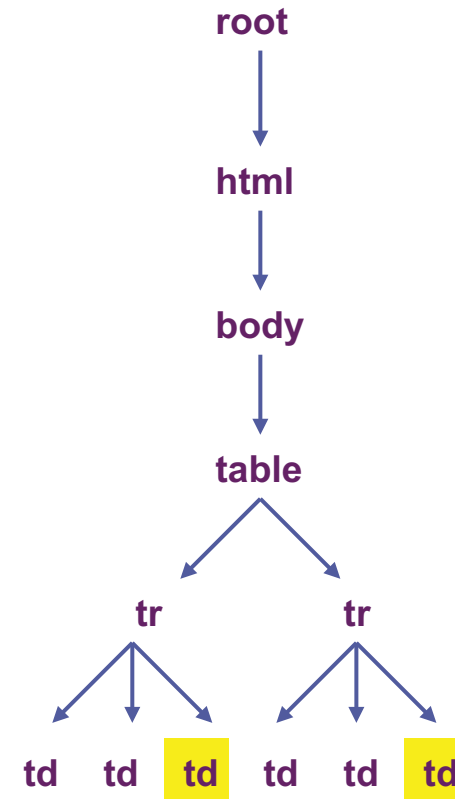
```
entry(X):- entry(Y), nextsibling(Y,X).
```

```
name(X) :- entry(E), firstchild(E, X), label[td](X).
```

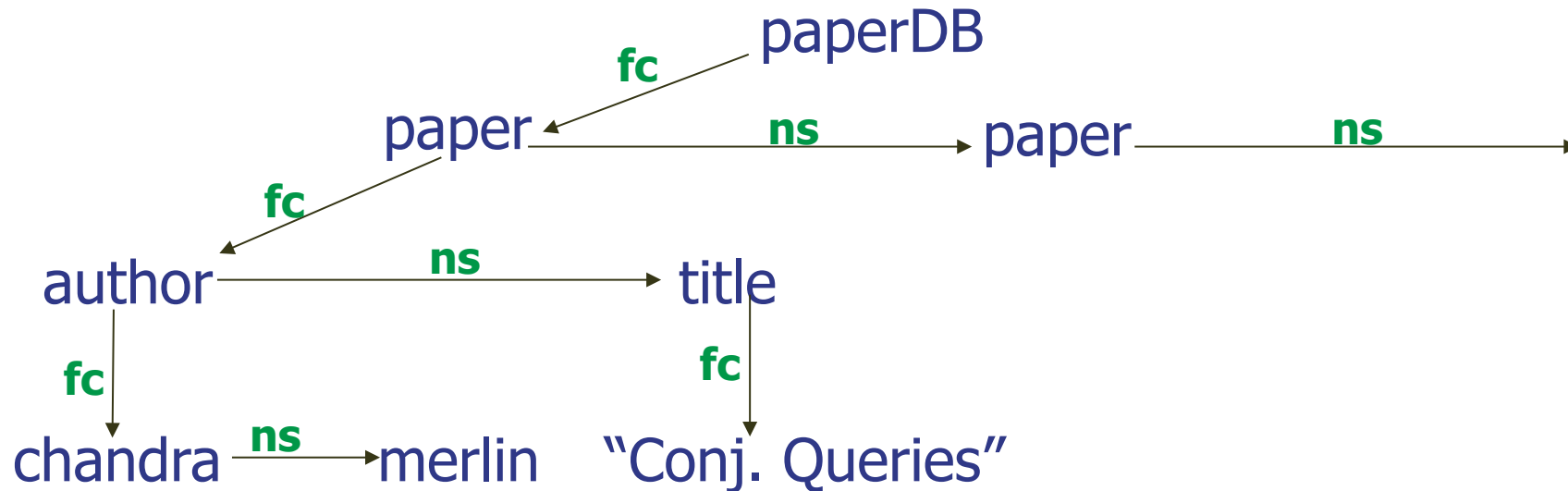
```
email(X) :- name(N), nextsibling(N, X), label[td](X).
```

```
phone(X) :- email(M), nextsibling(M, X), label[td](X).
```

```
<?xml version="1.0"?>
<peopledb>
<entry> <name>Georg Gottlob</name>
        <email>gottlob@dbai.tuwien.ac.at</email>
        <phone>18420</phone>
</entry>
<entry> <name>Christoph Koch</name>
        <email>koch@dbai.tuwien.ac.at</email>
        <phone>18449</phone>
</entry>
</peopledb>
```



# Monadic Datalog over XML



**Select titles of articles authored by Chandra and Merlin**

```
paper(X) ← root(R) & firstchild(R,X).  
paper(X) ← paper(Y) & nextsibling(Y,X).
```

```
output(X) ← paper(P) & firstchild(P,A) &  
firstchild(A,Z) & label[Chandra](Z) &  
nextsibling(Z,V) & label[Merlin](V) &  
nextsibling(A,T) & firstchild(T,X).
```

# How expressive is monadic Datalog?

It was known that over arbitrary structures:

- ◆ Monadic Datalog  $\subseteq \Pi_1$ -MSO
- ◆ Full Datalog = P (in presence of order)

**Theorem** [G. & Koch 2002]:

Over trees, monadic Datalog = MSO

A unary query is definable in MSO iff it is definable via a monadic datalog program.

# How complex is Monadic Datalog?

**Theorem** [G. & Koch 2002]:

Monadic Datalog over trees has  
combined complexity:  $O(|\text{data}| * |\text{query}|)$

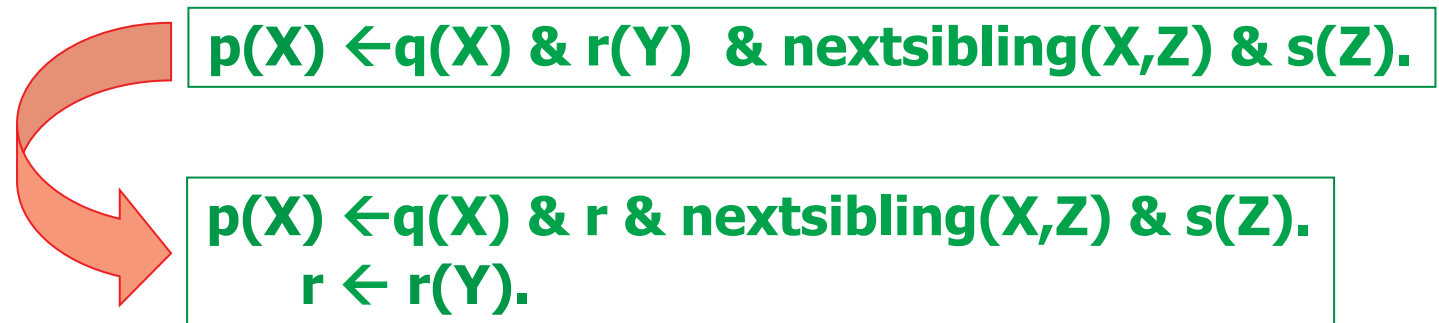
Query Complexity: P-complete and linear-time.



## Proof idea:

1.) Transform datalog program + input tree in linear time into a “ground” propositional logic program (programme Datalog instancié)

- Exploit functional dependencies:  
 $\text{nextsibling}(X,Y)$  has only a linear number of ground instances:  $\text{nextsibling}(n_i,n_j)$ , etc.
- Decouple independent atoms of rule bodies



2.) Execute ground program in linear time by using well-known algorithms: [Beeri&Bernstein][Dowling&Gallier] [Minoux]

**Logic heaven**

**DB theory heaven**

**DB programming  
heaven**

**Application design  
heaven**



**MSO**

||

**Monadic Datalog**

∩

**Elog**











∩

**Lixto Visual Wrapper**

Lixto Wrapper Generator - Extracted\_sources

Program Document

Extraction program: H:\lixto\reexamples\ebay\ebay.xml Active example document: file://h:/lixto/reexamples/ebay/ebay1.html

	<a href="#">NEC LAPTOP/NOTEBOOK PII-233,96MB,13", CD</a> <i>Buy It Now</i>	\$399.00	-	Jul-13 10:14
	<a href="#">NEC Versa S/500 486 DX 12 MB RAM 340 MB HD</a> 	\$20.60	7	Jul-17 09:57
	<a href="#">NEC READY 120LT NOTEBOOK COMPUTER</a> 	\$350.00	-	Jul-17 09:44
	<a href="#">NEC 200C W/IBM ETHERJET CARD!</a>  <i>Buy It Now</i>	\$50.00	-	Jul-17 09:30
	<a href="#">IBM 760XL CD, FLOPPY AND MORE!</a>  <i>Buy It Now</i>	\$99.00	-	Jul-17 09:12
	<a href="#">Qty 10 NEC Versa LX PII 300Mhz/128M warranty</a>	\$4400.00	-	Jul-17 09:04

For more items in this category, click these pages:  
= 1 = [2](#) [3](#) [4](#) [5](#) [\(next page\)](#)

Close

item description and link to detailpage

# of bids

date

one record

next page link

price info

# ELOG

[Baumgartner, Flesca, G. VLDB' 01]

## Examples of Special predicates:

**subelem**(S,X,Path,...)

Xpath-like expression

**before**(X,Y,.....)

**after**(X,Y,...)

distance tolerance,etc.

**property**(X,Attribute, Op, Value.....)

**document**(URL,D)

**getdocumentFromHref**(X,D),  
**etc.**

**Additional features:** Stratified negation,  
string processing  
ontological concepts “**phonenum**(X)”  
ranges: **H**(S,X) :- **body**(.....)[1,5]  
object hierarchies

notebook all of eBay - includes all regions Search [more search options](#)

Results by: THUNDERST

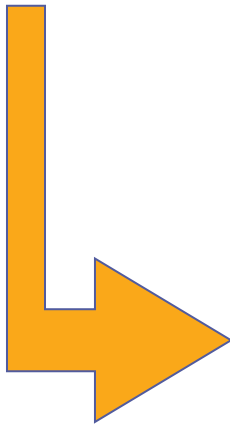
Search titles and descriptions (to find more items!) [Search Completed Items](#)

[eBay official time](#) 02:25:2

2150 items found for "notebook". Showing items 1 to 50.

Sort: Items ending first

Item#	Item	Price	Bids	Ends PDT
409449118	<a href="#">98 Degrees - Notebook - New</a> 🏠	\$2.99	-	in 19
413171469	<a href="#">Notebook - Compaq Presario 1207</a>	AU \$730.00	6	Aug-21 0
409454540	<a href="#">Compaq Armada Notebook P-100 Win 95</a>	\$107.50	8	Aug-21 0
409456450	<a href="#">THE NOTEBOOK NICHOLAS SPARKS HARDCOVER</a> 🏠	\$5.50	2	Aug-21 0



```
<?xml version="1.0" encoding="UTF-8"?>
<document>
  <record>
    <number>409449118</number>
    <item>98 Degrees - Notebook - New</item>
    <picture/>
    <price>2.99</price>
    <currency>$</currency>
    <bids>-</bids>
  </record>
  <record>
    <number>413171469</number>
    <item>Notebook - Compaq Presario 1207</item>
    <price>730.00</price>
    <currency>AU $</currency>
  </record>
  [...]

```

# ELOG Program for eBay pages

```
tableseq(S, X) ← document("www.ebay.com/", S), subseq(S, (.body, []), (.table, []), (.table, []), X),
                 before(S, X, (.table, [(elementtext, item,)], 1, 1, -, -), after(S, X, (.hr, []), 1, 1, -, -)

record(S, X) ← tableseq(_, S), subelem(S, .table, X)

itemnum(S, X) ← record(_, S), subelem(S, *.td, X), notbefore(S, X, (.td, []), maxint)

itemdes(S, X) ← record(_, S), subelem(S, (*.td *.content, [(a, , 0)], X)

price(S, X) ← record(_, S), subelem(S, (*.td, [(elementtext, Y, 1)]), X), valuta(Y)

bids(S, X) ← record(_, S), subelem(S, *.td, X), before(S, X, (.td, []), 1, 30, Y, _), price(S, Y)

currency(S, X) ← price(_, S), subtext(S, Y, X), valuta(Y)

pricewc(S, X) ← price(_, S), subtext(S, [0 - 9]+, X)
```

**Logic heaven**

**DB theory heaven**

**DB programming  
heaven**

**Application design  
heaven**



**MSO**

**||**

**Monadic Datalog**

**∩**

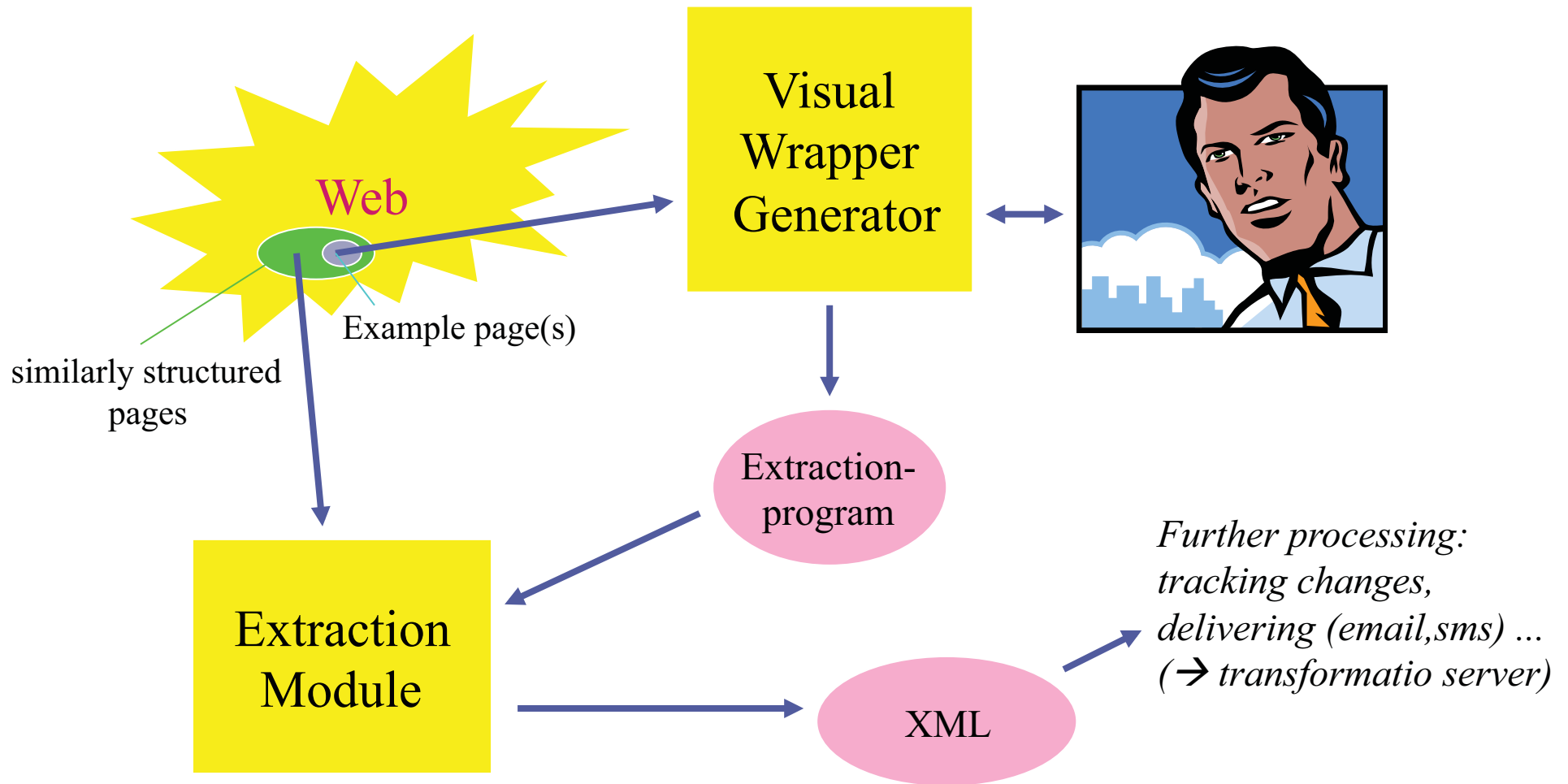
**Elog**

**∩**

**Lixto Visual Wrapper**

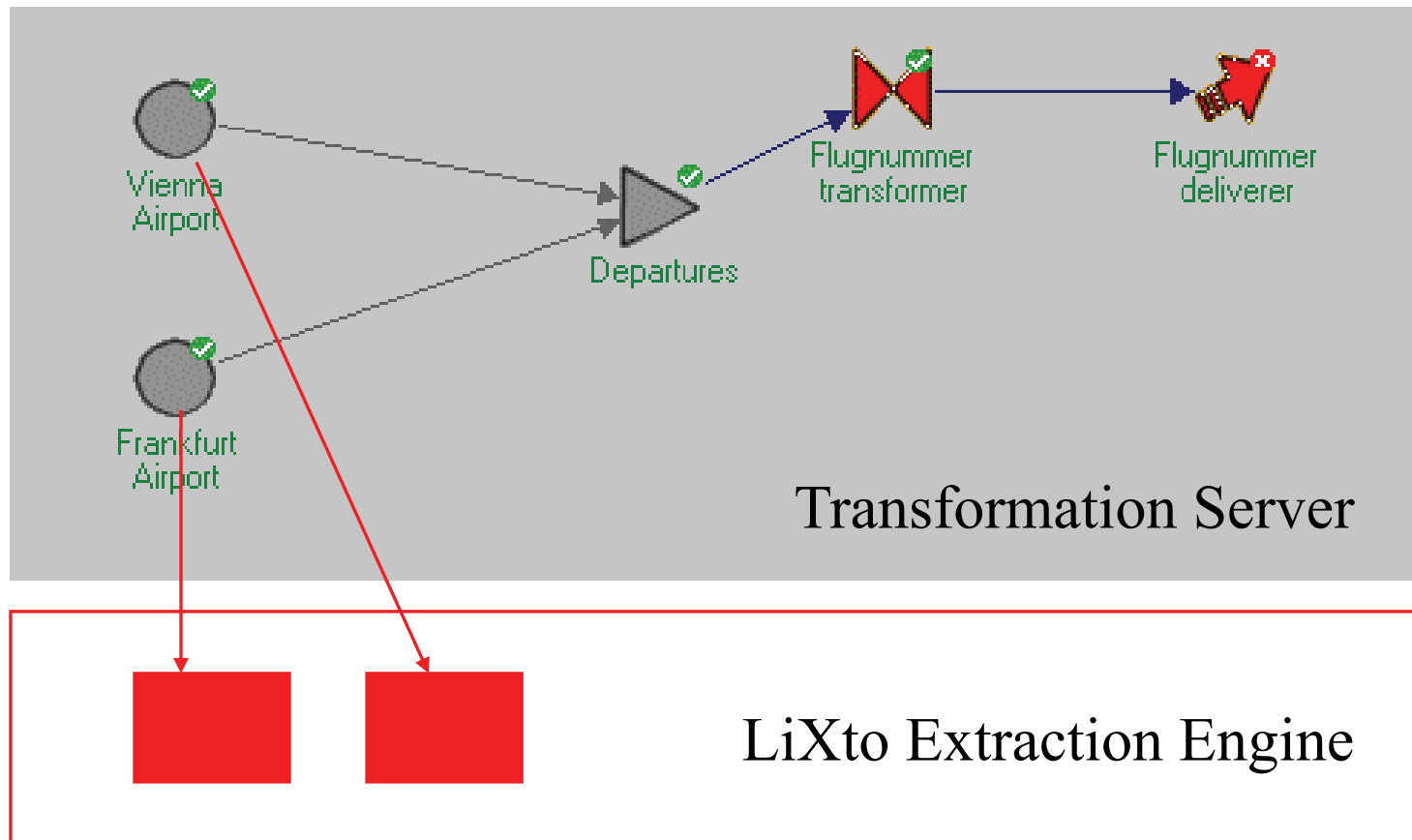
(outil: suite logicielle)

# Lixto Visual Wrapper Architecture





# Product Architecture



# **SHORT DEMO**

# Talk Outline

---

- Motivation: need of information extraction
- Logical foundations of information extraction
- The Lixto Visual Wrapper
- **The Diadem Project: Fully automatic data extraction**

## Need for Automatic Extraction Technology (2)

**All search engine providers need it! Many work on it.**

**Keywords:** → Vertical search,  
→ object search,  
→ semantic search.

**Raghu Ramakrishnan**, Yahoo!, March 2009:

*"no one really has done this successfully at scale yet"*

**Alon Halevy**, Google, Feb. 2009: *"Current technologies are not good enough yet to provide what search engines really need. [...] any successful approach would probably need a combination of knowledge and learning."*



**ERC Advanced Grant  
Research proposal (Part B1)**

# **Domain-centric Intelligent Automated Data Extraction Methodology**

## **DIADEM**

**Principal Investigator:** Georg Gottlob

**Host Institution:** University of Oxford

**Full title:** Domain-centric Intelligent Automated Data Extraction Methodology

**Short name:** DIADEM

**Duration in months:** 60

*Gottlob*



**ERC Advanced Grant  
Research proposal (Part B1)**

## **Domain-centric Intelligent Automated Data Extraction Methodology**

### **DIADEM**

**Principal Investigator:** Georg Gottlob

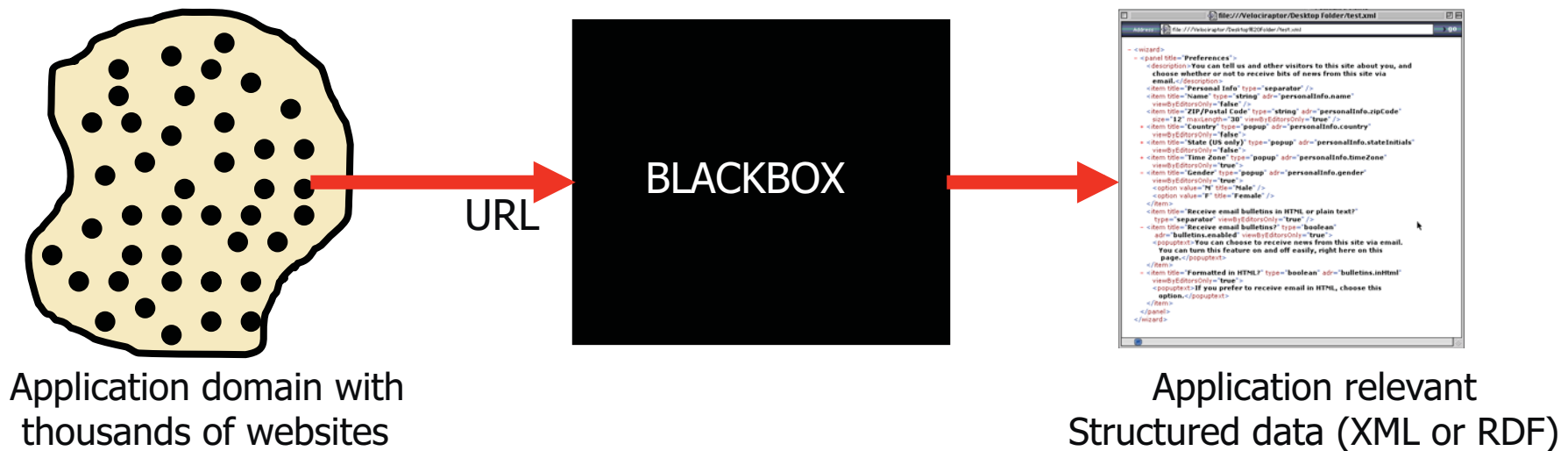
**Host Institution:** University of Oxford

**Full title:** Domain-centric Intelligent Automated Data Extraction Methodology

**Short name:** DIADEM

**Duration in months:** 60

# The Blackbox we are constructing



To achieve this, we combine a host of annotators with a new knowledge-based approach.

# How to achieve it?

Combine existing and new “low level” annotators with “high level” AI and reasoning.



## Property Search

UK

International

I'm interested in

- Buying
- Renting
- New developments

Maximum price

GBP ▾ Any ▾

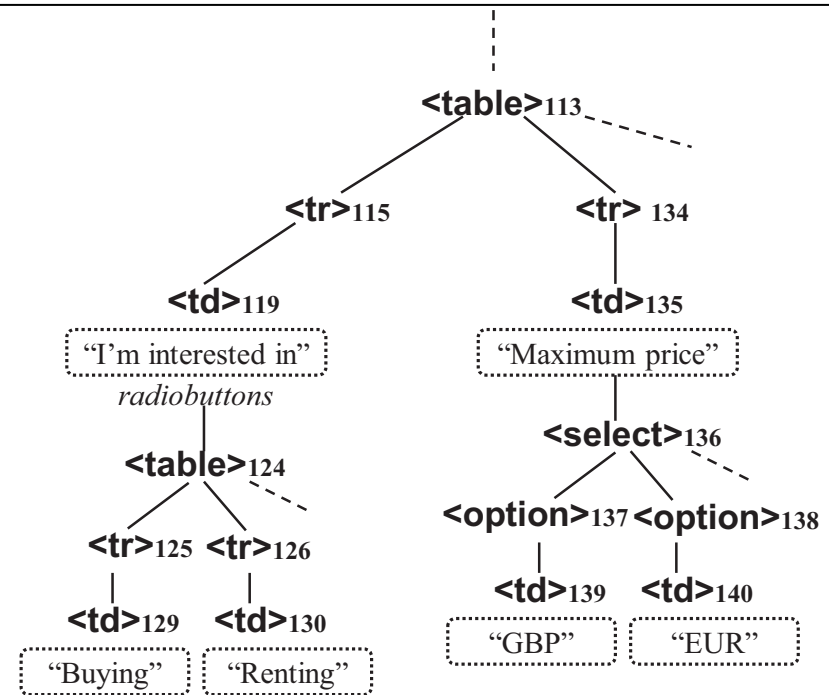
Minimum bedrooms

Any ▾

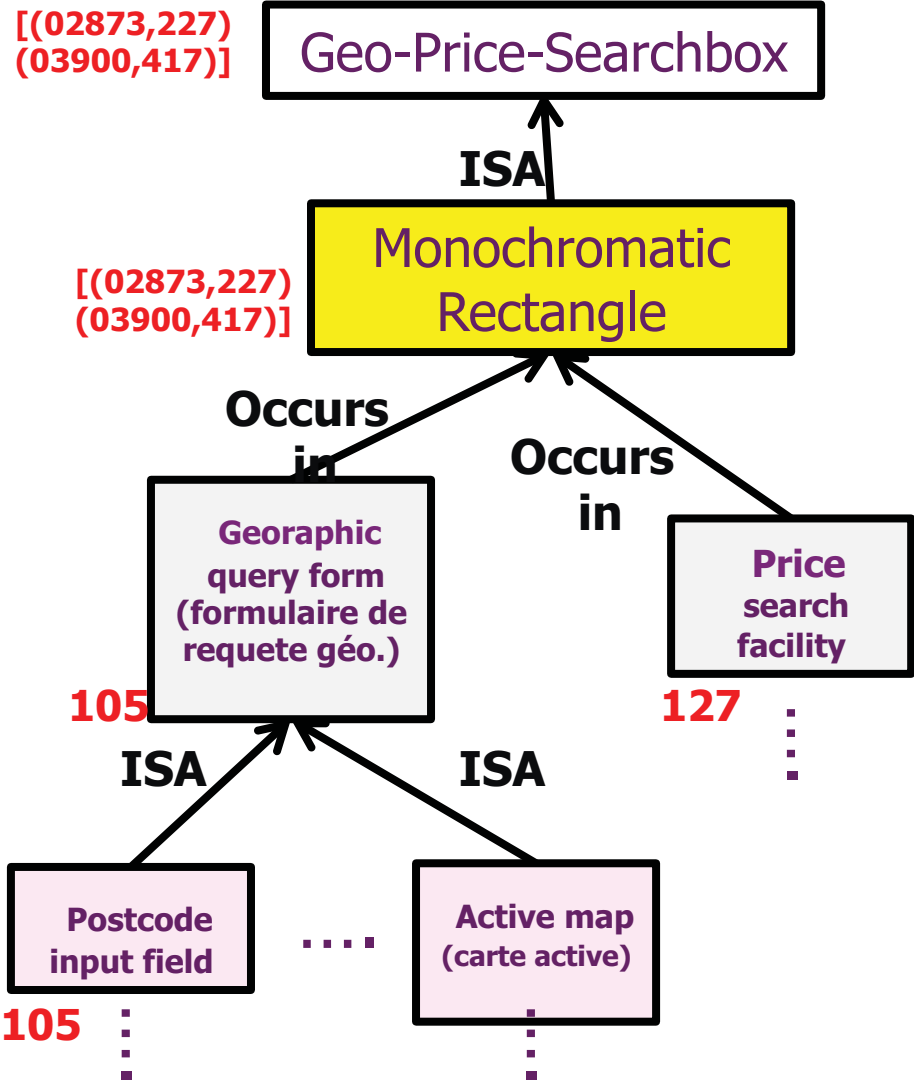
Area

e.g Surrey, SW4, Chelsea

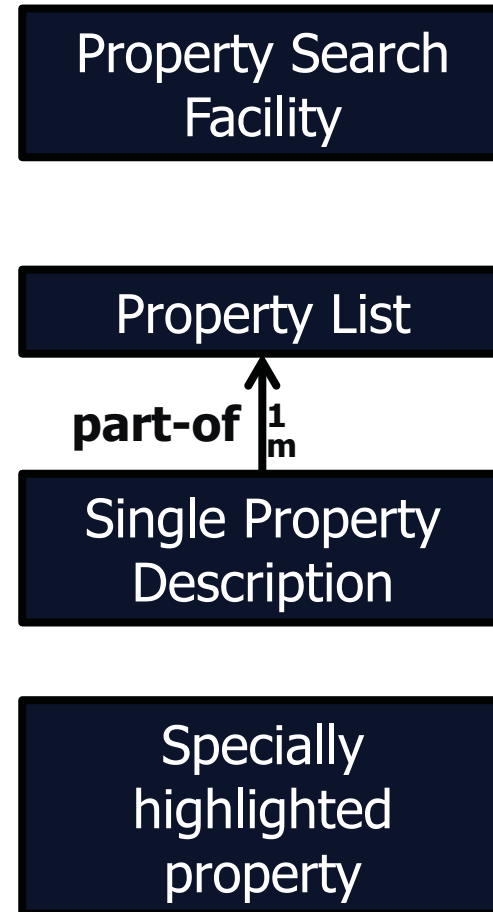
GO



# Bottom-up (low-level) annotation



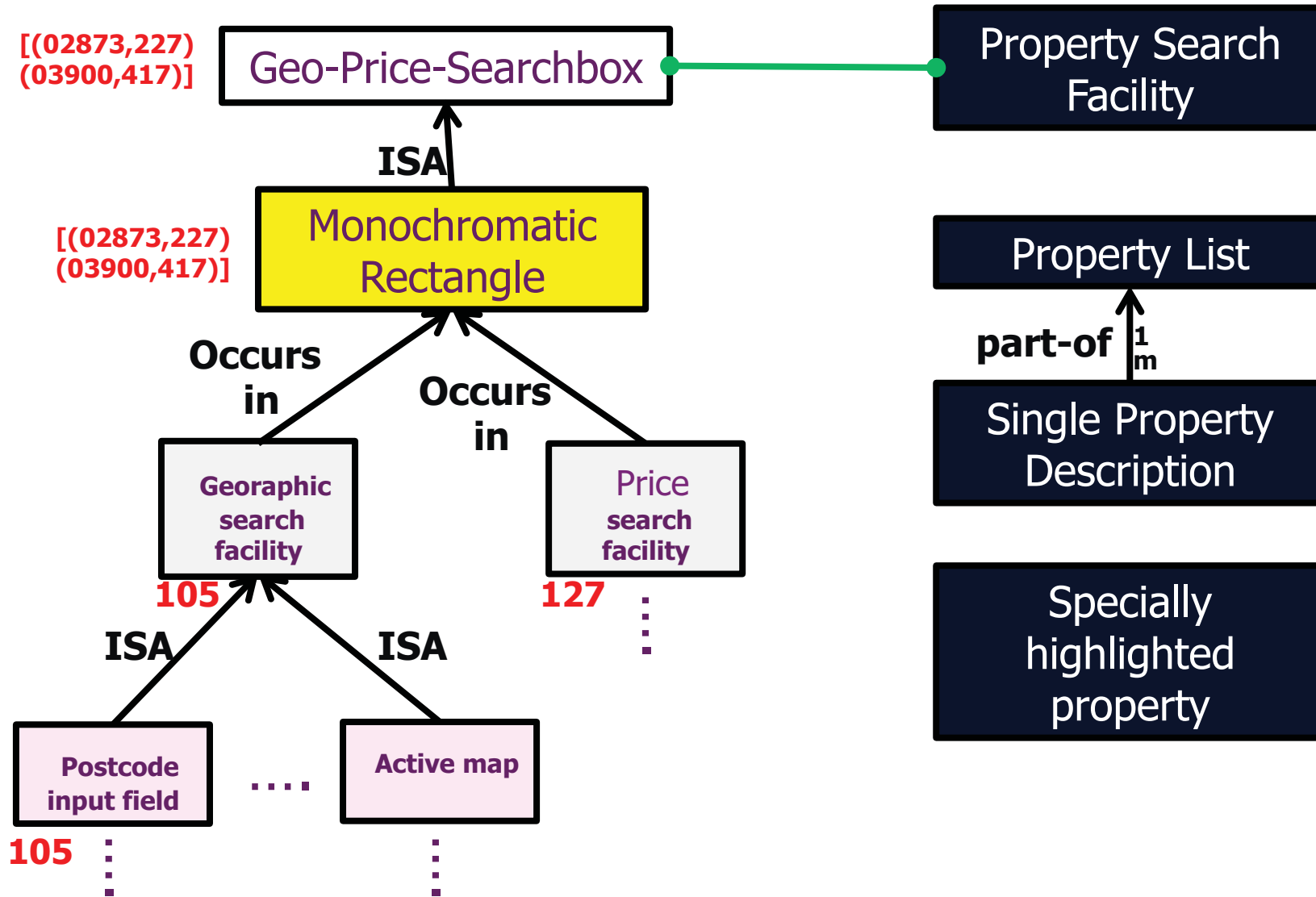
# Top-down reasoning



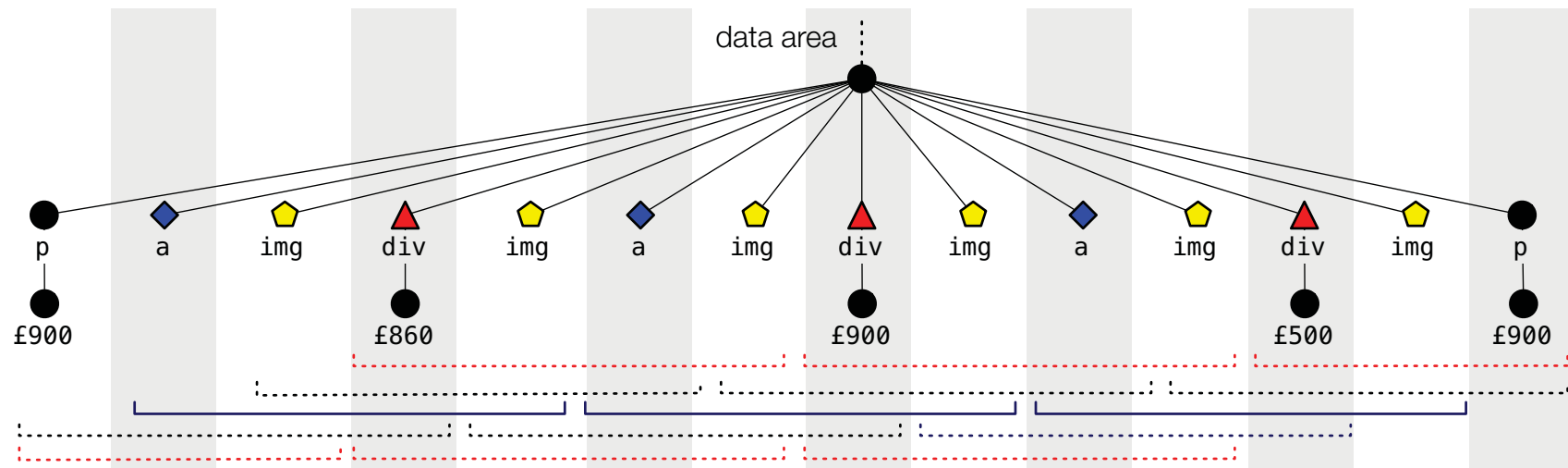
Bottom-up processing

Top-down reasoning

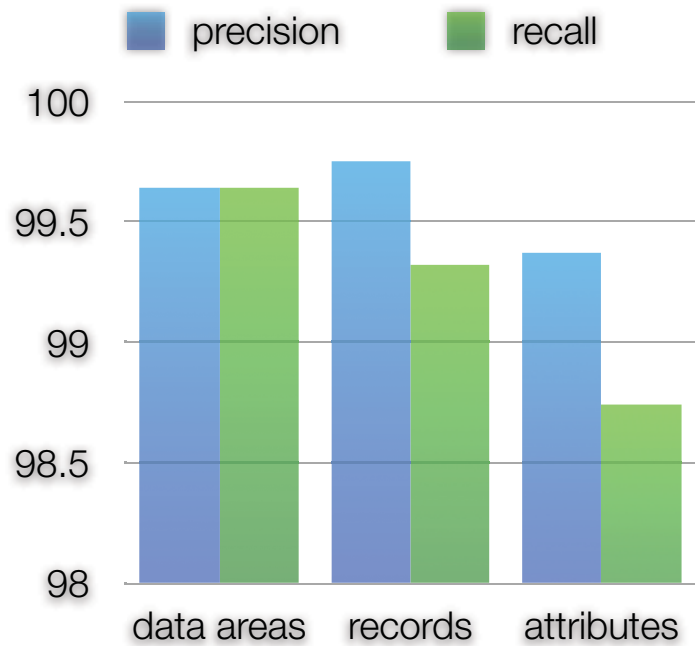
### Phenomenology



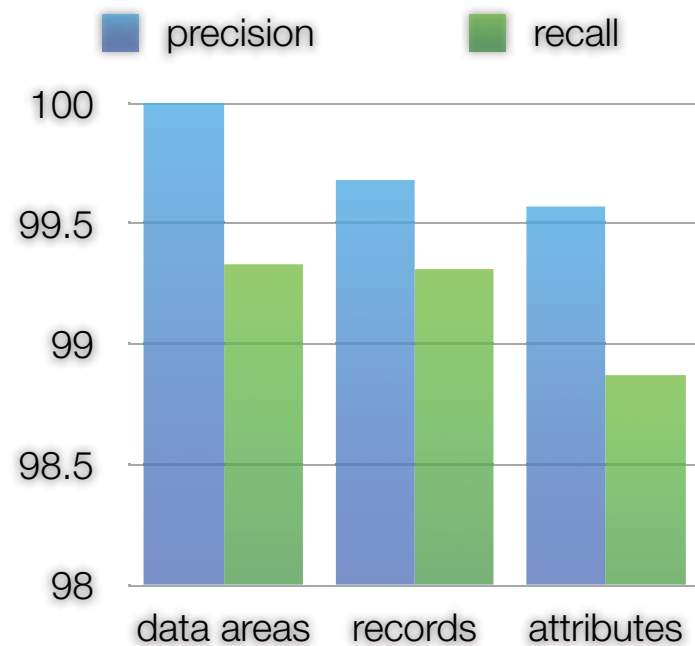
# Phenomenological Record Segmentation



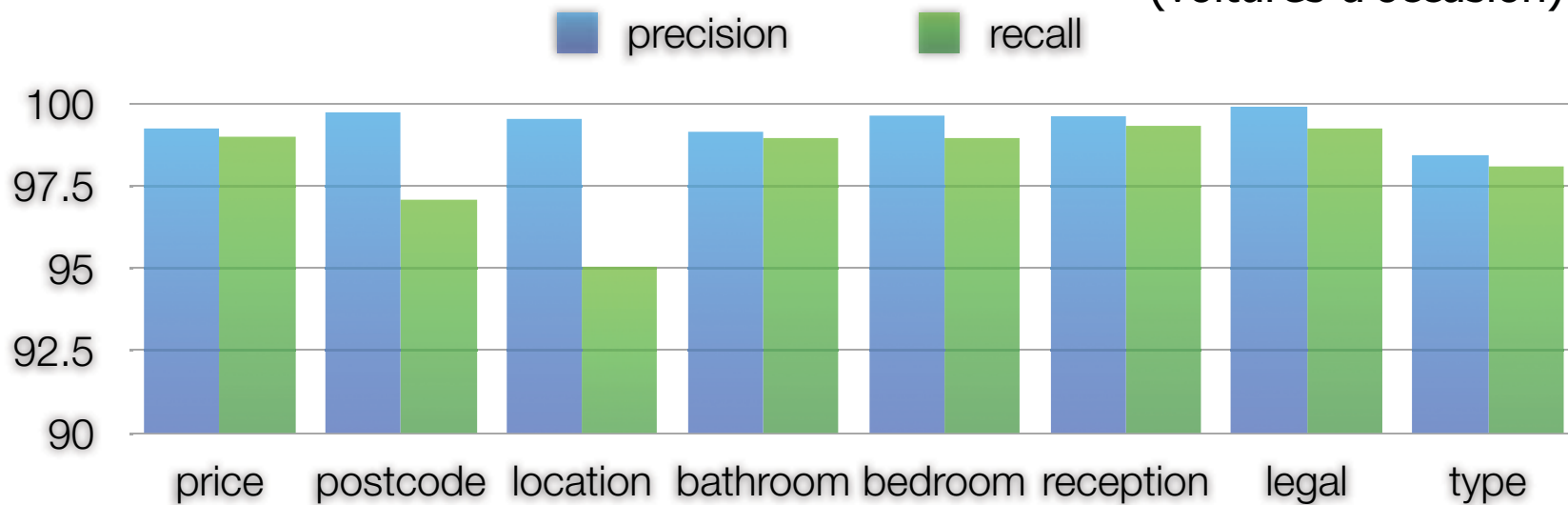
- set of uniform, non-overlapping records
- maximise sequence of evenly segmented (same distance pivot)
- minimise irregularity of records



**Real Estate**  
(100 pages)



**Used Car**  
(100 pages)  
(voitures d'occasion)



3

# Form Patterns Example

The image shows a car search interface with several highlighted form patterns:

- Price (£) (?):** A range slider with input boxes for minimum (0) and maximum (50000) values, and an "Exclude 'Call For Price'" checkbox.
- Age (years):** A range slider with input boxes for minimum (0) and maximum (10+) values.
- Mileage (miles):** A range slider with input boxes for minimum (0) and maximum (100K+) values.
- Number of owners:** A radio button selection list with options: Any, 1 previous owner (0), up to 2 previous owners (0), up to 3 previous owners (0), more than 3 previous owners, and Unknown (0).
- PROPERTY SEARCH:** A sub-form with a "Price" section containing "No min" and "No max" dropdowns, and a "No. of beds" dropdown set to "All".

- Small set of ubiquitous patterns
  - ranges, dates, options, etc.
- Ontology by instantiation

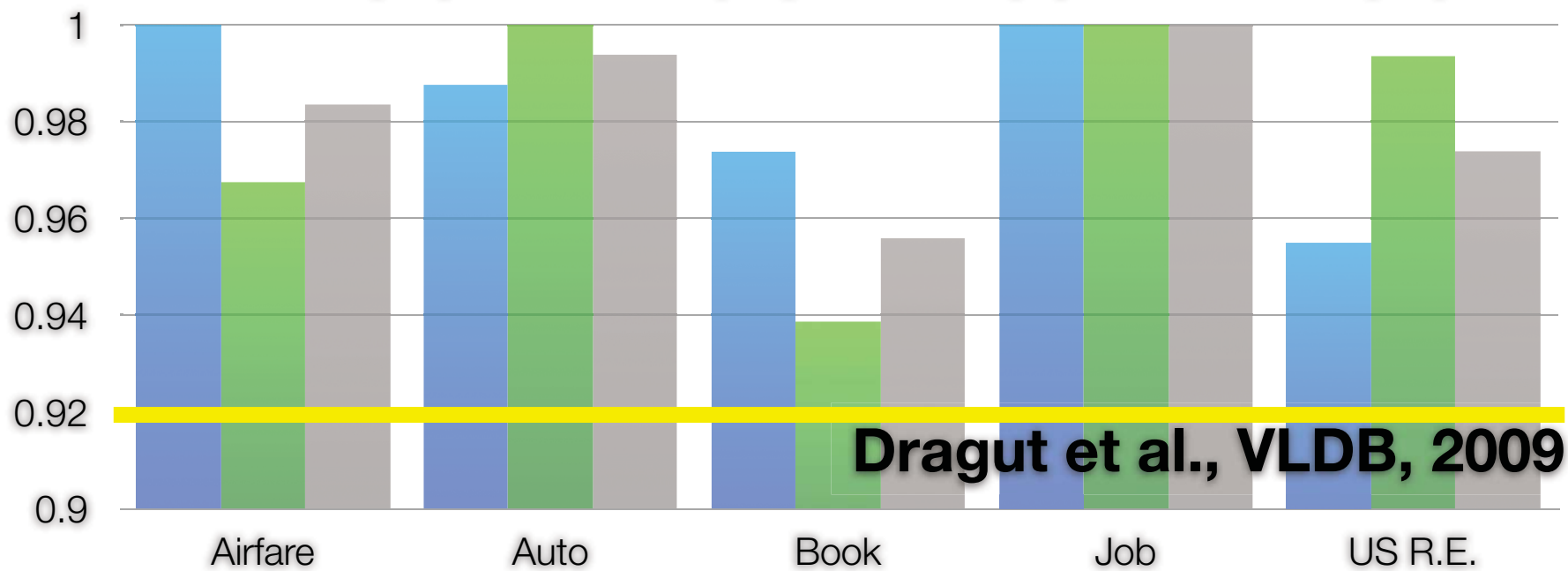
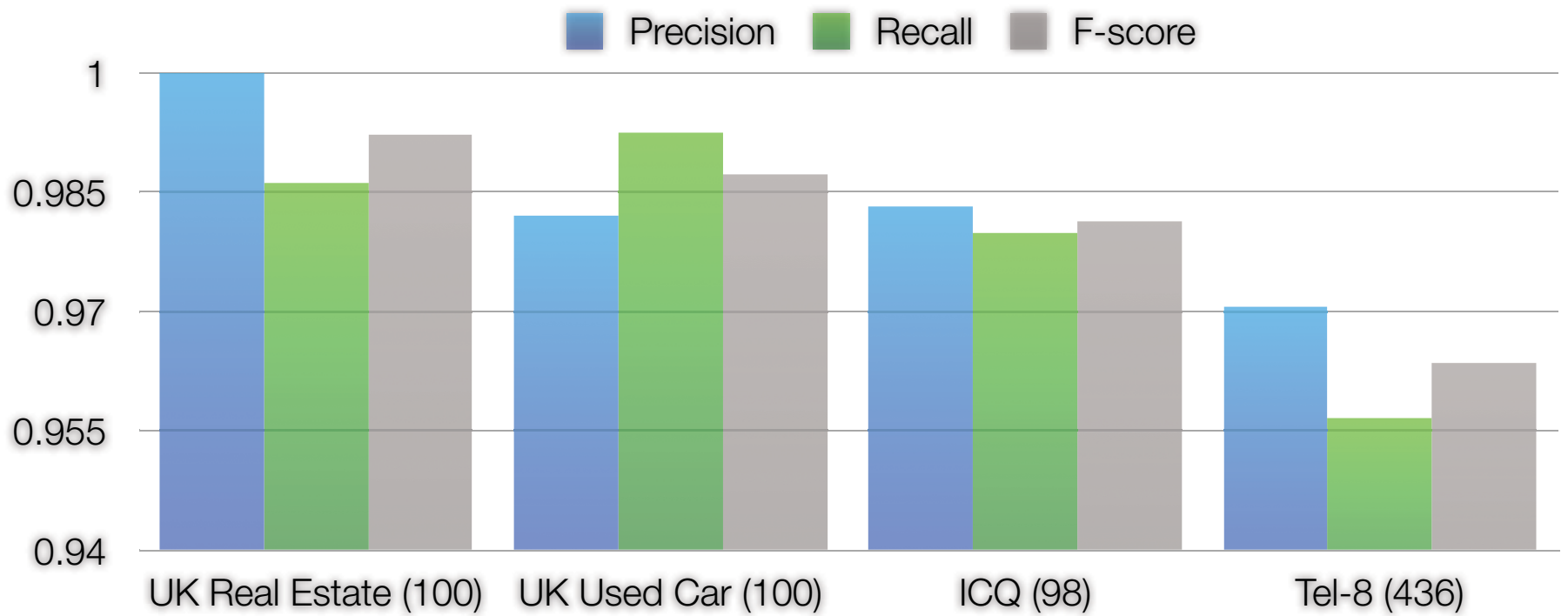
# OPAL-TL Example

- Price range
  - two successive fields in the same group
  - at least one “price” type
  - range connector in between

The image shows a 'PROPERTY SEARCH' form with several input fields. A red rectangular box highlights the 'Price' section, which contains two dropdown menus: 'No min' and 'No max', separated by the word 'to'. Other visible fields include 'Sales' (selected) and 'Lettings' (unselected) radio buttons, a 'Town / City' dropdown menu set to 'All', and a 'No. of beds' dropdown menu set to 'All'. There are also radio buttons for 'List view' (selected) and 'Map view' (unselected), and a 'Search' button at the bottom right.

$\text{concept}\langle C_M \rangle(N_2) \Leftarrow \text{child}(N_1, G), \text{child}(N_2, G), \text{follows}(N_2, N_1),$   
 $\text{concept}\langle C \rangle(N_1), N_2 @ \text{range\_connector}\{e, d\}, \neg(A_1 \prec A, N_2 @ A_1 \{d\})$





Short Demo [diadem-3min43.m4v](#)