# My Adventures with Datalog

## Walking the Thin Line Between Theory and Practice

Georg Gottlob

University of Oxford & TU Wien

# Cultural Background



H. Hahn, M. Schlick, O. Neurath, R. Carnap, K. Gödel



Ludwig Wittgenstein

Vienna was a hotspot of logic in the early 20$^{th}$ century.

This was stopped by the rise of the Nazis in the late 1930ies.

**Goal**: without any pretention, help reestablishing a tradition of logic in Austria: **Logic in CS – our Mission**
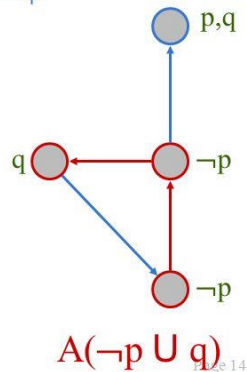
**Computer Science is the continuation of logic by other means.**

# New logics for specific applications have been designed

Examples:

### Temporal logic (CTL)

◆ Propositional logic:
  – p, true, false, $\varphi \wedge \psi$, $\varphi \vee \psi$, $\neg\varphi$
◆ One step in the future
  – One path: $EX\varphi$
  – All paths: $AX\varphi$
◆ Multiple steps in the future
  – One path: $E(\varphi \cup \psi)$
  – All paths: $A(\varphi \cup \psi)$
◆ Backwards variants
  – $AX\triangle\varphi$, $A\triangle(\varphi \cup \psi)$, etc.

p,q

q   ¬p

¬p

$A(\neg p \cup q)$



### Datalog: Logics for Big Data

**Deductive Database**

• Explicit information
• Rules that enable inferences based on the stored data

parent(x,y)

| | |
|---|---|
| Alice | Nancy |
| Alice | Joyce |
| Joyce | Lois |
| Lois | Mark |
| Lois | Andy |
| Joyce | Ruth |

Datalog program

anc(x,y) :- parent(x,y)
anc(x,y) :- anc(x,z), parent(z,y)

head          body
recursions

$\forall$ x,y   (anc(x,y) ← parent(x,y))
$\forall$ x,y,z (anc(x,y) ← anc(x,z), parent(z,y))
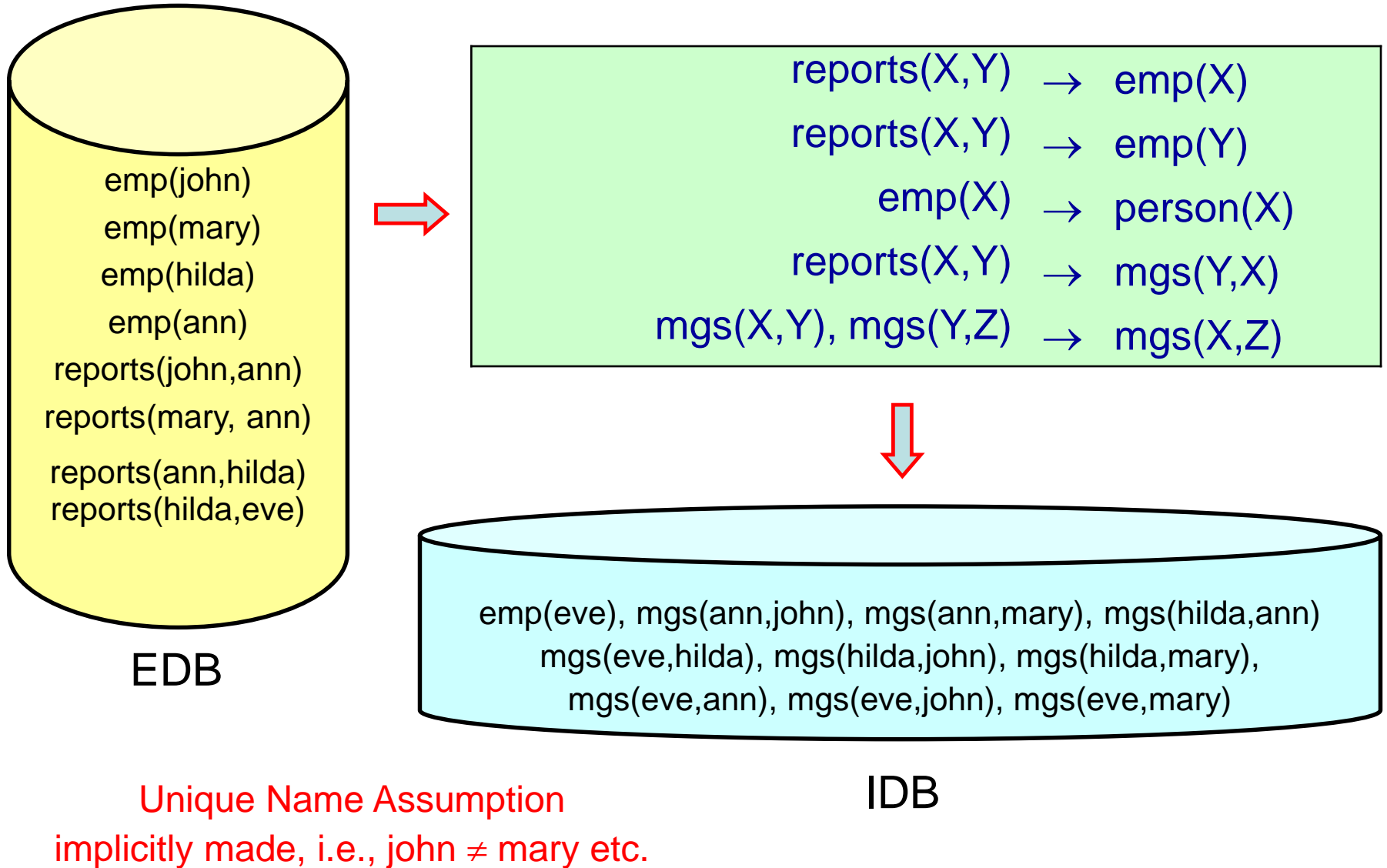


Reasoning with these logics is <u>decidable</u>, their exact complexity is known,
and they can be directly <u>implemented</u> as computer programs.
→ Suitable for industrial applications.

## Make Logic Work!

Many variants and many other logics/algorithms (e.g. **SAT-Solving, OWL**) have been designed.

# Datalog - An Example

EDB (yellow cylinder):
- emp(john)
- emp(mary)
- emp(hilda)
- emp(ann)
- reports(john,ann)
- reports(mary, ann)
- reports(ann,hilda)
- reports(hilda,eve)

**EDB**

Rules (green box):
$$reports(X,Y) \rightarrow emp(X)$$
$$reports(X,Y) \rightarrow emp(Y)$$
$$emp(X) \rightarrow person(X)$$
$$reports(X,Y) \rightarrow mgs(Y,X)$$
$$mgs(X,Y), mgs(Y,Z) \rightarrow mgs(X,Z)$$

IDB (blue cylinder):
emp(eve), mgs(ann,john), mgs(ann,mary), mgs(hilda,ann)
mgs(eve,hilda), mgs(hilda,john), mgs(hilda,mary),
mgs(eve,ann), mgs(eve,john), mgs(eve,mary)

**IDB**

Unique Name Assumption
implicitly made, i.e., john $\neq$ mary etc.

# Logic in Industry…

… is dramatically gaining momentum, and more jobs in LogiCS are becoming available.

**SIEMENS**
**Austria**

**LogicBlox**
**Predictix**

Automatic configuration of:

- Railway safety systems, telecommunication systems, etc.
- Model-based methods, answer-set programming, SAT solving.

- Based on **Datalog**
- Software for business applications
- about 100 employees including 20-30 computational logicians
- $40m investment
- Recently sold to INFOR

# Datalog Startup

**semmle >**

- Based on **Datalog**
- Software engineering, particularly **code analysis** and **code querying**
- More than fifty employees
- Founded 2006 by **Oege de Moor** as spin-out of Oxford University (total investment $21M)
- Now headquarter in San Francisco and offices in Oxford, New York City, and Copenhagen.
- Eager to hire people whose skills include:
  - Database theory
  - Declarative/logic programming (particularly **Datalog**)
  - Formal logic, lattice theory, abstract interpretation

*Recently sold to Github/Microsoft*

# Logic in the Industry



**diffblue**
*We write tests for you. Automatically*

Founded 2016 by **Daniel Kroening** as spin-out of Oxford University.

Logic-based **verification tools** and **software test automation**.

Funding: $22M by four investors;

Eagerly looking for LogiCS PhDs.



## AMAZON'S AUTOMATED REASONING GROUP IS HIRING!

Ready to solve big problems, and work on innovative projects that will affect millions of customers worldwide? Amazon Web Services is looking for intelligent, customer obsessed, research driven graduate students within the focused areas of computer-aided verification and programming languages for full-time and internship positions within our Automated Reasoning Group.

**Amazon's automated reasoning team is a global initiative:**

Seattle - WA | **New York** - NY | **Washington** - DC | **Bangalore** - India | **Dresden** - Germany | **London** - UK

Automated reasoning drives security and operational excellence across AWS!
We are hiring Formal Verification Engineers to work in some of the following areas:

Code Verification - Theorem Proving - Programming Languages -
- Computer-aided verification - Satisfiability Modulo Theories - and more!

# My own Experience as Founder: Co-founded 4 Companies

| **Projects at TU Wien** | **Projects at TU Wien** | **Projects at Oxford U.** | **Projects at Oxford U.** |
|---|---|---|---|
| *CD Laboratory f. Expert Syst.* (CDG 1989-96) | *Wittgenstein Award* (FWF 1998-2004) | *DIADEM* (ERC 2009-2014) | *Schema Mappings* (EPSRC 2007-2010) |
| Nicola Leone: *Query System f. Disjunctive DBs* (FWF 1996-2000) | *Inductive Learning for Web Data Extraction* (FWF Transl. 2005-8) | *James Martin Grant* (J.M. School 2014-15) | *VADA: Value-Added Data* (EPSRC 2015-2020) |
| | | *ExtraLytics* (ERC 2014-16) | |

**DLVSYSTEM**

**Main Founder: Nicola Leone**

1998

liXto

2001

Wrapidity

2015

DeepReason.ai

2017

McKinsey & Company

Periscope® By McKinsey

2013

Meltwater Outside Insight

2016

1. From Disjunctive Datalog Research to DLVSYSTEM

# 1. From Disjunctive Datalog Research to DLV-System



**edge(X,Y) → edge(Y,X)**
**edge(X,Y) → node(X)**
**edge(X,Y) → node(Y)**
**node(X) → color(X, red) ∨ color(X, blue) ∨ color(X, green)**
**color(X, red) ∧ color(X, blue) → false**
**color(X, red) ∧ color(X, green) → false**
**color(X, green) ∧ color(X, blue) → false**
**edge(X,Y) ∧ color(X,C) ∧ color(Y,C) → false**

# 1. From Disjunctive Datalog Research to DLV-System



**edge(X,Y) → edge(Y,X)**
**edge(X,Y) → node(X)**
**edge(X,Y) → node(Y)**
**node(X) → color(X, red) ∨ color(X, blue) ∨ color(X, green)**
**color(X, red) ∧ color(X, blue) → false**
**color(X, red) ∧ color(X, green) → false**
**color(X, green) ∧ color(X, blue) → false**
**edge(X,Y) ∧ color(X,C) ∧ color(Y,C) → false**

**edge(X,Y) → edge(Y,X)**
**edge(X,Y) → node(X)**
**edge(X,Y) → node(Y)**
**node(X) → color(X, red) ∨ color(X, blue) ∨ color(X, green)**
**color(X, red) ∧ color(X, blue) → false**
**color(X, red) ∧ color(X, green) → false**
**color(X, green) ∧ color(X, blue) → false**
**edge(X,Y) ∧ color(X,C) ∧ color(Y,C) → false**

Example of reasoning question:
Are there necessarily two nodes of every color?

# 1. From Disjunctive Datalog Research to DLV-System



edge(X,Y) → edge(Y,X)
edge(X,Y) → node(X)
edge(X,Y) → node(Y)
node(X) → color(X, red) ∨ color(X, blue) ∨ color(X, green)
color(X, red) ∧ color(X, blue) → false
color(X, red) ∧ color(X, green) → false
color(X, green) ∧ color(X, blue) → false
edge(X,Y) ∧ color(X,C) ∧ color(Y,C) → false

## Disjunctive Datalog

THOMAS EITER
University of Giessen
GEORG GOTTLOB
Technical University of Vienna
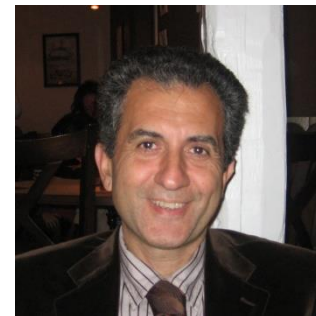and
HEIKKI MANNILA
University of Helsinki

We consider disjunctive Datalog, a powerful database query language based on disjunctive gic programming. Briefly, disjunctive Datalog is a variant of Datalog where disjunctions may appear in the rule heads; advanced versions also allow for negation in the bodies, which can be handled according to a semantics for negation in disjunctive logic programming. In particular, we investigate three different semantics for disjunctive Datalog: the minimal model semantics, the perfect model semantics, and the stable model semantics. For each of these semantics, the expressive power and complexity are studied. We show that the possibility variants of these semantics express the same set of queries. In fact, they precisely capture the complexity class $\Sigma_2^p$. Thus, unless the Polynomial Hierarchy collapses, disjunctive Datalog is more expressive than normal logic programming with negation. These results are not only of theoretical interest; we demonstrate that problems relevant in practice such as computing the optimal tour value in the Traveling Salesman Problem and eigenvector

ACM Transact. on Database Systems 22(3): 364-418 (1997)

Theorem:
Disjunctive Datalog reasoning is $\Pi_2^p$–complete and can express every property in $\Pi_2^p$.

**D L V S Y S T E M** S.R.L.
SPIN-OFF OF UNIVERSITY OF CALABRIA

**DLV**

DLV is an artificial intelligence system based on disjunctive logic programming, which offers front-ends to several advanced KR formalisms.

**DLV$^{DB}$**

DLV$^{DB}$ is an extension of the DLV system designed both to handle input and output data distributed on several databases.

**ASPIDE**

ASPIDE is a Integrated Development Environment for Answer Set Programming supporting the entire life-cycle of ASP development.

**JDLV**

JDLV is a new programming framework blending DLV with Java programming.

## Example

| Program Π | | | Query $Q$ |
|---|---|---|---|
| $a \leftarrow not\ b$ | $b \leftarrow not\ a$ | % either $a$ or $b$ | $a, b, c$ |
| $c \leftarrow a$ | $c \leftarrow b$ | | |
| $d \leftarrow not\ e$ | $e \leftarrow not\ d$ | % either $d$ or $e$ | |
| $\leftarrow a, c, d$ | | % added after step 1 | |
| $\leftarrow a, c, e$ | | % added after step 2 | |
| $\leftarrow b, c, d$ | | % added after step 3 | |
| $\leftarrow b, c, e$ | | % added after step 4 | |

### Execution

| Step | Stable model | Underestimate | Overestimate |
|---|---|---|---|
| 0 | | $\emptyset$ | $\{a, b, c\}$ |
| 1 | $\{a, c, d\}$ | $\emptyset$ | $\{a, c\}$ |
| 2 | $\{a, c, e\}$ | $\emptyset$ | $\{a, c\}$ |
| 3 | $\{b, c, d\}$ | $\emptyset$ | $\{c\}$ |
| 4 | $\{b, c, e\}$ | $\emptyset$ | $\{c\}$ |

## 2. From Monadic Datalog Research to Lixto

# Web Data Extraction



| ref-code | postcode | bedrooms | bathrooms | available | price |
|----------|----------|----------|-----------|-----------|----------|
| 33453 | OX2 6AR | 3 | 2 | 15/10/2013 | £1280 pcm |
| 33433 | OX4 7DG | 2 | 1 | 18/04/2013 | £995 pcm |

# Function  f: HTML Parse tree ➔ Subtrees

# Leaves of subtrees are among leaves of orig. tree

# 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>

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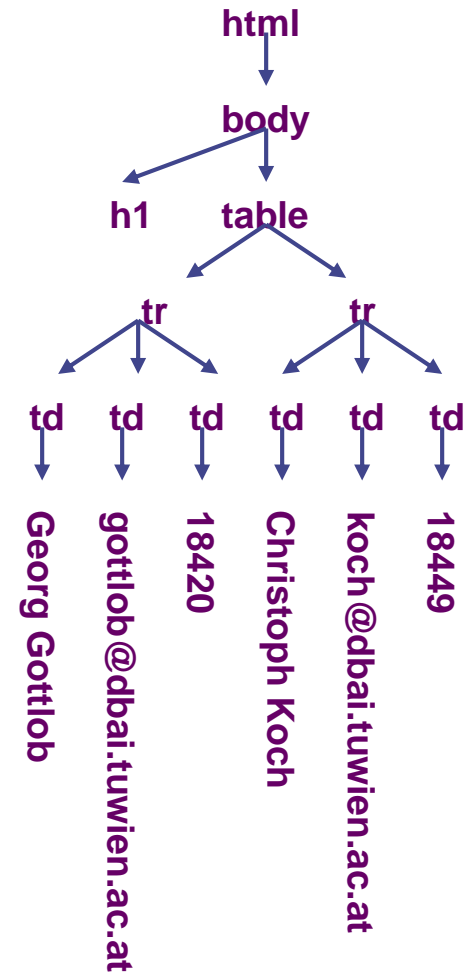
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN">

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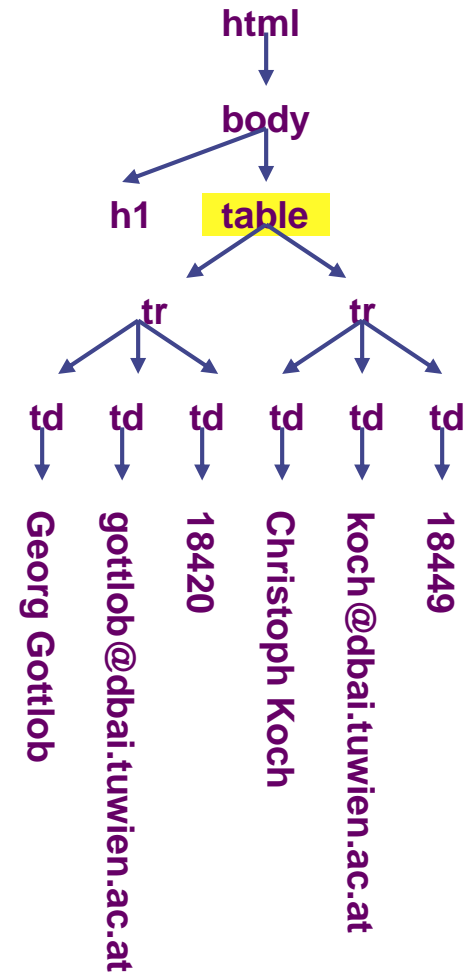
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN">

<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 |

# 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).

root

html

body

table

tr        tr

td   td   td   td   td   td

# 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).**

root

html

body

table

tr        tr

td  td  td   td  td  td

# 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).**

root

html

body

table

tr        tr

td  td  td   td  td  td

# 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).

root

html

body

table

tr          tr

td  td  td   td  td  td

# Monadic Datalog as a Wrapping Language
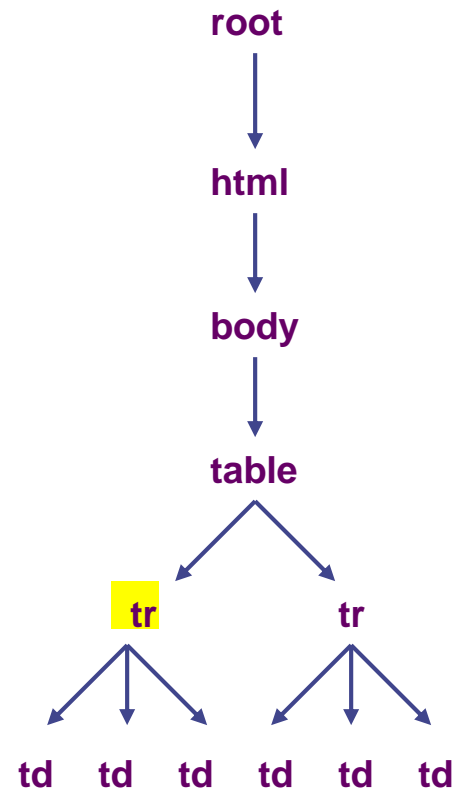
entry(X) :- root(R), firstchild(R,U), label[html](U),
              firstchild(U,V), label[body](V),
              firstchild(V,W),label[table](W),
              firstchild(W,X), label[tr](X).

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).

root

html

body

table

tr          tr

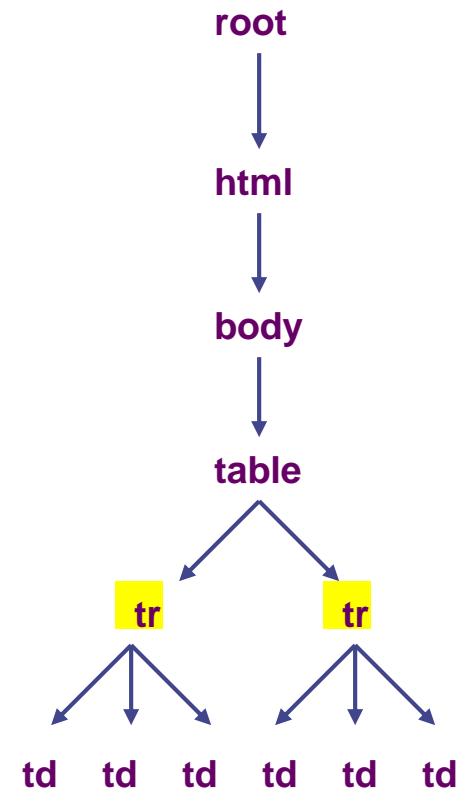td  **td**  td    td  **td**  td

# Monadic Datalog as a Wrapping Language

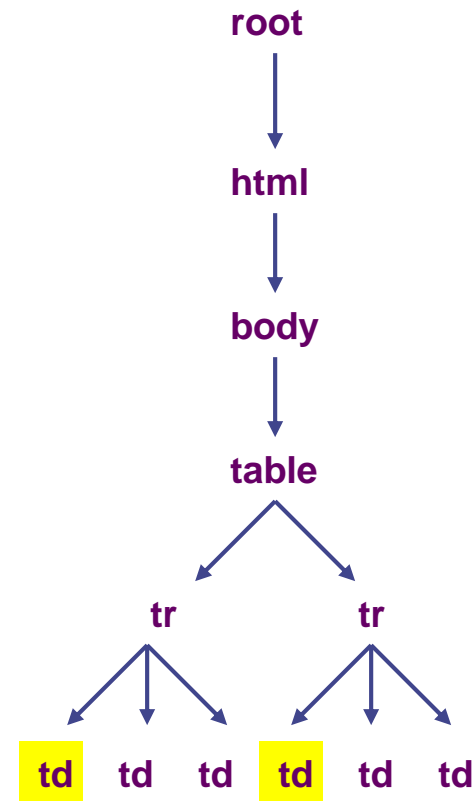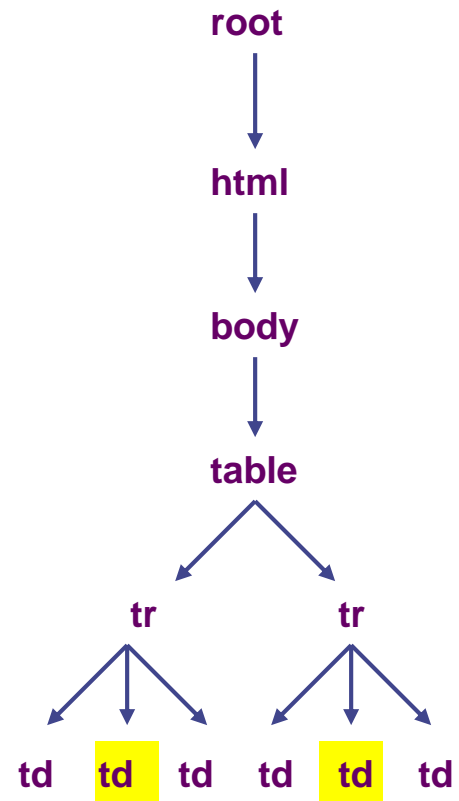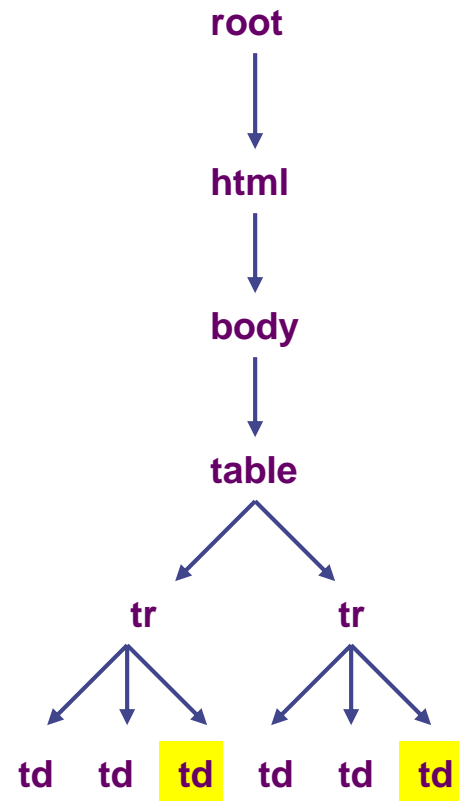entry(X) :- root(R), firstchild(R,U), label[html](U),
                   firstchild(U,V), label[body](V),
                   firstchild(V,W),label[table](W),
                   firstchild(W,X), label[tr](X).

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).

# Theoretical Results

**Theorem1: Monadic Datalog over trees has combined complexity: $O(|data| * |query|)$**

Monadic Datalog and the Expressive Power of
Languages for Web Information Extraction

GEORG GOTTLOB and CHRISTOPH KOCH
Technische Universität Wien, Austria

Research on information extraction from Web pages (wrapping) has seen much activity recently (particularly systems implementations), but little work has been done on formally studying the expressiveness of the formalisms proposed or on the theoretical foundations of wrapping. In this paper, we first study monadic datalog over trees as a wrapping language. We show that this simple language is equivalent to monadic second order logic (MSO) in its ability to specify wrappers. We believe that MSO has the right expressiveness required for Web information extraction and propose MSO as a yardstick for evaluating and comparing wrappers. Along the way, several other results on the complexity of query evaluation and query containment for monadic datalog over trees are established, and a simple normal form for this language is presented. Using the above results, we subsequently study the kernel fragment Elog⁻ of the Elog wrapping language used in the Lixto system (a visual wrapper generator). Curiously, Elog⁻ exactly captures MSO, yet is easier to use. Indeed, programs in this language can be entirely visually specified.

## 1. INTRODUCTION

The Web wrapping problem, i.e., the problem of extracting structured information from HTML documents, is one of high practical importance and has spurred a great

**Theorem 2 : Over trees, Monadic Datalog = MSO**

MSO= Monadic second order logic.

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

# ELOG Program for eBay pages

$tableseq(S, X) \leftarrow document(\text{"www.ebay.com/"}, S), subsq(S, (.body, []), (.table, []), (.table, []), X),$
$before(S, X, (.table, [(elementtext, item, )]), 1, 1, \_, \_), after(S, X, (.hr, []), 1, 1, \_, \_)$

$record(S, X) \leftarrow tableseq(\_, S), subelem(S, .table, X)$

$itemnum(S, X) \leftarrow record(\_, S), subelem(S, \star.td, X), notbefore(S, X, (.td, []), maxint)$

$itemdes(S, X) \leftarrow record(\_, S), subelem(S, (\star.td. \star .content, [(a, , 0)], X)$

$price(S, X) \leftarrow record(\_, S), subelem(S, (\star.td, [(elementtext, Y, 1)]), X), valuta(Y)$

$bids(S, X) \leftarrow record(\_, S), subelem(S, \star.td, X), before(S, X, (.td, []), 1, 30, Y, \_), price(S, Y)$

$currency(S, X) \leftarrow price(\_, S), subtext(S, Y, X), valuta(Y)$

$pricewc(S, X) \leftarrow price(\_, S), subtext(S, [0 - 9]^{+}, X)$

# Lixto Visual Developer (VD)



Mozilla Web Browser

Navigation Steps

Extraction Configuration

Extern·Extraction Cluster

amazon
web services™

HydraNT Server

Transformation Server

liXto
THE WEB INTELLIGENCE COMPANY

Data Center

Intern·Extraction Cluster

## Search

| | |
|---|---|
| Product Group: | All Product Groups ▼ |
| EAN/UPC: | |
| Product Name: | |
| Week: | Any Week ▼ |

Search

## Prices by Different Sources

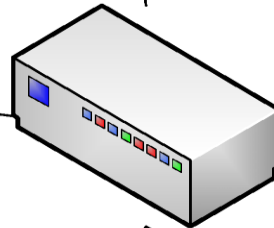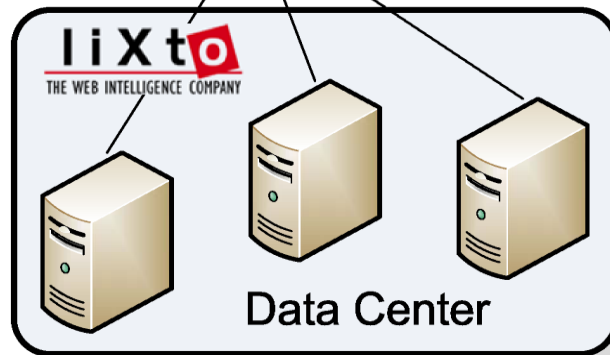| Week ▲ | EAN/UPC | Product Name | You | Source1 | Source2 | Source3 | Source4 | Source5 | Source6 | Source7 | Min | Max | Avg | Meetbeat | Diff | Diff [%] |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 38 | 4045827397809 | Product #432 | | | | 1070€ | 799€ | 798€ | | | 798€ | 1070€ | 889€ | | | |
| 38 | 4045827401063 | Product #437 | 799€ | | | | 749€ | 798€ | | | 749€ | 799€ | 782€ | 0% | +50€ | +7% |
| 38 | 4045827404835 | Product #438 | | | 890€ | | | 780€ | | | 780€ | 890€ | 835€ | | | |
| 38 | 4045827406969 | Product #442 | 999€ | | | | | 1065€ | | | 999€ | 1065€ | 1032€ | 100% | -66€ | -6% |
| 38 | 4045827407010 | Product #444 | 899€ | | | | | 1300€ | | | 899€ | 1300€ | 1100€ | 100% | -401€ | -31% |
| 38 | 4045827409533 | Product #445 | | | 1128€ | | | 1445€ | | | 1128€ | 1445€ | 1287€ | | | |
| 38 | 4045827417095 | Product #446 | | | | 832€ | 849€ | 885€ | 899€ | | 832€ | 899€ | 866€ | | | |
| 38 | 4045827432715 | Product #460 | 899€ | | 914€ | 832€ | | 1298€ | | | 832€ | 1298€ | 986€ | 67% | +67€ | +8% |
| 38 | 4045827439509 | Product #462 | 999€ | | | 796€ | 999€ | 998€ | | | 796€ | 999€ | 948€ | 0% | +203€ | +26% |
| 38 | 4045827449553 | Product #477 | | | | 1426€ | 929€ | | 1099€ | | 929€ | 1426€ | 1151€ | | | |
| 38 | 4045827453604 | Product #509 | | 594€ | 474€ | | | 700€ | | | 474€ | 700€ | 589€ | | | |
| 38 | 4045827466666 | Product #540 | | | | | | 690€ | 1199€ | | 690€ | 1199€ | 945€ | | | |
| 38 | 4045827466697 | Product #545 | | 713€ | | | 1599€ | | | | 713€ | 1599€ | 1156€ | | | |
| 38 | 4045827476481 | Product #547 | 1499€ | | | | | 1499€ | | | 1499€ | 1499€ | 1499€ | 0% | +0€ | +0% |
| 38 | 4045827492764 | Product #581 | 699€ | | | | | 683€ | | | 683€ | 699€ | 691€ | 0% | +16€ | +2% |
| 38 | 4045827492771 | Product #583 | | | | 689€ | | 998€ | | | 689€ | 998€ | 844€ | | | |
| 38 | 4045827492788 | Product #585 | 799€ | | | | | 775€ | | | 775€ | 799€ | 787€ | 0% | +24€ | +3% |
| 38 | 4045827493150 | Product #587 | 649€ | | 831€ | | | 785€ | | | 649€ | 831€ | 755€ | 100% | -136€ | -17% |
| 38 | 4045827506690 | Product #629 | 499€ | | | | | 498€ | | | 498€ | 499€ | 499€ | 0% | +1€ | +0% |
| 38 | 4045827511090 | Product #642 | | | | 594€ | 599€ | | | | 594€ | 599€ | 597€ | | | |

row(s) 1 - 20 of 467 ▼   Next ▶

# *August 2013:*

## McKinsey & Company acquires Lixto

**Periscope, a McKinsey Solution, has enhanced its suite of revenue management solutions with McKinsey & Company's acquisition of Lixto. McKinsey …**

Periscope, a McKinsey Solution, has enhanced its suite of revenue management solutions with McKinsey & Company's acquisition of Lixto.

McKinsey & Company, a global management consulting firm, has announced the acquisition of Lixto Software, a Vienna-based SaaS solutions company in the field of online competitive intelligence.

This acquisition will enhance **Periscope**, an integrated suite of solutions within McKinsey Solutions designed to deliver sustainable ROS improvement through better pricing, promotions, assortment and performance management.

Lixto will continue to offer services for automated web data extraction, data mapping and matching, providing frequent and reliable competitive online data on pricing, promotions, assortment range, terms/conditions and product attributes, delivering a competitive advantage to a variety of clients in Retail, Travel & Hospitality, High Tech, Consumer Products, Distribution, and Manufacturing industries. These capabilities will support clients with market intelligence and actionable insights, powered by Periscope Performance Vision market analytics and Periscope Price Advisor dynamic pricing solutions.

3. From Lixto over DIADEM to Wrapidity

This [Lixto] is great technology, but …
… too many people needed for generating and
maintaining the wrappers!

From the tourism, real estate, and retail domains, we understood independently that *automated* wrapper-generation would be useful.

But how can this´be achieved?

# Need for Fully Automated Extraction Technology

Example: Real Estate UK > 15000 sites

  many not covered by aggregators

  list of all agencies easy to get (source discovery)

  but: manual or semi-automatic wrapping too expensive

    wrapper construction

    testing

    tracking changes

No existing tool or methodology could do it fully automatically

DIADEM

a €2.5M project
started in April '10

Application domain with thousands of websites

Application-relevant Structured data (XML or RDF)

# Rough Idea: Knowledge via Datalog Rules

Use "expert" rules that analyze Web pages and interact with them

- **Ontological rules**  (how do entities relate to each other)

  - a flat is a real-estate property

  - a house is a real-estate property

  - a real-estate property has a number of rooms

  - a price consists of a number and a currency

- **Phemomenological rules**  (how do entities manifest themselves on the Web?)

  - the text chunk closest to an input field is with high prob. its descriptor.

  - each sales item is described in a "convex"  (usual. rectangular) region.

- **Site exploration rules:**

  - before filling a field try to leave it empty

  - rules for handling next-page links

- **Other types of rules**

# DIADEM PROJECT

- DIADEM lab at Oxford University

| 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | VADA → |

*spin-out start-up*
*Wrapidity*

Funding so far > $7M

# Research Papers (examples)

## DIADEM: Thousands of Websites to a Single Database*

Tim Furche, Georg Gottlob, Giovanni Grasso, Xiaonan Guo,
Giorgio Orsi, Christian Schallhart, Cheng Wang

Department of Computer Science, Oxford University, Wolfson Building, Parks Road, Oxford OX1 3QD
firstname.lastname@cs.ox.ac.uk

### ABSTRACT

The web is overflowing with implicitly structured data, spread over hundreds of thousands of sites, hidden deep behind search forms, or siloed in marketplaces, only accessible as HTML. Automatic extraction of structured data at the scale of thousands of websites has long proven elusive, despite its central role in the "web of data".

Through an extensive evaluation spanning over 10000 web sites from multiple application domains, we show that automatic, yet accurate full-site extraction is no longer a distant dream. DIADEM is the first automatic full-site extraction system that is able to extract structured data from different domains at very high accuracy. It combines automated *exploration* of websites, *identification* of relevant data, and *induction* of exhaustive wrappers. Automating these components is the first challenge. DIADEM overcomes this challenge by combining phenomenological and ontological knowledge. Integrating these components is the second challenge. DIADEM overcomes this challenge through a self-adaptive network of relational transducers that produces effective wrappers for a wide variety of websites.

Our extensive and publicly available evaluation shows that, for more than 90% of sites from three domains, DIADEM obtains an effective wrapper that extracts all relevant data with 97% average precision. DIADEM also tolerates noisy entity recognisers, and its components individually outperform comparable approaches.

### Categories and Subject Descriptors

H.3.5 [**Information Storage and Retrieval**]: On-line Information Services—*Web-based services*

### General Terms

Languages, Experimentation

### Keywords

data extraction, deep web, wrapper induction

### 1. INTRODUCTION

The web has become the largest repository of structured data. For the US alone, the number of online shopping sites with $10k+ revenue is estimated in excess of a hundred thousand [30], with a long-tail of several hundred thousand smaller shops. A significant portion of data is only available in this long-tail [11].

This data is mostly available in HTML pages, designed for humans. The automatic, yet accurate extraction of the structured data underlying such pages is a long standing challenge [5]. Semi-supervised data extraction approaches, such as [2, 12], have been investigated extensively, but require users to supervise the induction by navigating each site and identifying relevant data. In contrast, **automatic full-site extraction** (AFE) operates *automatically* with no per-site supervision, navigates to all relevant data on the *full site*, yet *extracts* highly *structured data*.

Automatic full-site extraction involves three primary sub-problems, namely site exploration (with form understanding and filling), record and attribute identification, and wrapper induction. Each of these sub-problems is a significant challenge by itself, but worse, these problems have in the past been tackled in isolation with few exceptions. Successful applications of AFE have been limited to narrow settings with simple structure, such as title and body extraction for news articles [39] or search engine results (ViNTs [44]). For extracting highly structured data, these approaches are unsuitable. Furthermore, most of these approaches fail on modern sites. For example, ViNTs [44] full-site extraction identifies records (of title and body only) with only 83%-88% accuracy (Section 5), even when supervised by selecting negative and positive examples of result pages for each site.

This lack of integrated, full-site extraction approaches has significantly reduced the impact of data extraction—with some commercial applications such as Google Products moving away from extraction towards purely curated data collection. Yet, the need has only increased, in particular with the rise of big data analytics for competitive price intelligence or intelligent supply chain management. In many of these applications, the acquisition of accurate, large-scale data, e.g., about competitor's products, current deals, or

**VLDB 2014**

---

SPECIAL ISSUE PAPER

## OXPATH: A language for scalable data extraction, automation, and crawling on the deep web

Tim Furche · Georg Gottlob · Giovanni Grasso ·
Christian Schallhart · Andrew Sellers

**Abstract** The evolution of the web has outpaced itself: A growing wealth of information and increasingly sophisticated interfaces necessitate automated processing, yet existing automation and data extraction technologies have been overwhelmed by this very growth. To address this trend, we identify four key requirements for web data extraction, automation, and (focused) web crawling: (1) interact with sophisticated web application interfaces, (2) precisely capture the relevant data to be extracted, (3) scale with the number of visited pages, and (4) readily embed into existing web technologies. We introduce OXPATH as an extension of XPATH for interacting with web applications and extracting data thus revealed—matching all the above requirements. OXPATH's page-at-a-time evaluation guarantees memory use independent of the number of visited pages, yet remains polynomial in time. We experimentally validate the theoretical complexity and demonstrate that OXPATH's resource consumption is dominated by page rendering in the underlying browser. With an extensive study of sublanguages and properties of OXPATH, we pinpoint the effect of specific features on evaluation performance. Our experiments show that OXPATH outperforms existing commercial and academic data extraction tools by a wide margin.

**Keywords** Web extraction · Crawling · Data extraction · Automation · XPath · DOM · AJAX · Web applications

### 1 Introduction

The dream that the wealth of information on the web is easily accessible to everyone is at the heart of the current evolution of the web. Due to the web's rapid growth, humans can no longer find all relevant data without automation. Indeed, many invaluable web services, such as Amazon, Facebook, or Pandora, already offer limited automation, focusing on filtering or recommendation. But in many cases, we cannot expect data providers to comply with yet another interface designed for automatic processing. Neither can we afford to wait another decade for publishers to implement these interfaces. Rather, data should be extracted from existing human-oriented user interfaces. This lessens the burden for providers, yet allows automated processing of everything accessible to human users, not just arbitrary fragments exposed by providers. This approach complements initiatives, such as Linked Open Data, which push providers toward publishing in open, interlinked formats.

For automation, data accessible to humans through existing interfaces must be transformed into structured data, for example, each gray span with class source on Google

T. Furche (✉) · G. Gottlob · G. Grasso · C. Schallhart · A. Sellers
Department of Computer Science, Oxford University,

**VLDB Journal 2013**

# DIADEM Architecture

| OPAL | AMBER | BERyL | OXPath |
|------|-------|-------|--------|
| Form filling & understanding | Object identification & alignment | Block analysis & object enrichment | Efficient extraction in the cloud |

**DATALOG± (implemented in DLV)**
Rule, exploration, control and integration language

New knowledge-based technology combining formalised knowledge with machine learning

**Application domain with thousands of websites**

**Application-relevant highly structured data**

Properties     Used cars 

Combines ML with rule-based Reasoning.

In production: Extracts from 200K websites in various application domains.

Successful on > 95% of test sites.

# Evaluation on 10k+ Sites

**10,493**
**Sites** from real-estate and used-car

**45**
**Node** Amazon EC2 cluster running 2.1 days

**92%**
**Effective wrappers** for more than 92% of sites on average

**97%**
**Precision** of extracted primary attributes

**100**
**Domain-dependent** concepts and relations

**20**
**Days** (one expert) to adjust system to a new domain

# Domains considered so far (since 2014)

- **Real estate UK**

- **Real estate US**

- **Used cars**

- **Consumer electronics**

- **Restaurant chains**

- **Restaurants in the 'Open Web'**

- **Jobs (from company Web sites)**

- **News**

- **Companies**

# Commercial Impact

ERC Advanced Grant DIADEM
+ ERC Proof of Concept Grant
EXTRALYTICS

Founded February 2015
operating initially in Oxford
now in London

2 possibilities:

- Build up company with large client portfolio

- Sell technology, IP & software to strategic partner

M Inbox (84,103) - gg ×   W Commercialization ×   Meltwater acquires ×   MTS Meltwater Acquires ×   iTWire - Meltwater ×   cso Meltwater acquir

← → C  ⓘ martechseries.com/analytics/meltwater-acquires-ai-wrapidity-to-add-ai-to-media-monitoring-capabilities/

⠿ Apps   Drucken - Proxy List F   One.com Webhosting   The best & free FTP S   Gerät(HY002017)   Lock Picking   G youtube zbigniew ryb   Imagir

**AI/ML**

# Meltwater Acquires Wrapidity to Add AI Capabilities into Media Intelligence Platform

By Sudipto Ghosh 🐦

Posted on February 21, 2017

Meltwater, the leading B2B data analytics company, has acquired London-based web data extraction company Wrapidity for an undisclosed amount. The AI-startup that spun out of Oxford University in 2015 will be a separate entity in Meltwater's existing platform.  By beefing up its "media intelligence" platform, Meltwater will now offer AI-powered automation tools for data analytics and media monitoring from unstructured web-based content.

In the era of specialized AI for MarTech, Wrapidity offers tailor-made solutions to content- specific problems arising in image recognition, Natural Language Processing, and machine learning. By acquiring Wrapidity, Meltwater will be able to automate its data extraction processes to reach out to a wide range of online customers based on accurate analytics of historical and real-time data. Meltwater is expected to further improve Wrapidity's AI capabilities for content discovery and data asset management, enabling marketers to interrogate data for diverse purposes, including sales enablement, social media monitoring and so on.
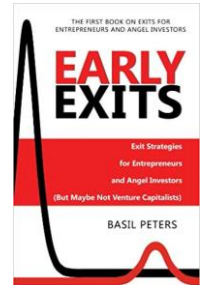
Meltwater empowers marketing teams

# Remarks to Prospective Academic Founders

**Question:** Big exit (8+ years)  vs  moderate exit (2-3 years)  ?

Big exit: *expected value* lower: Risks (see Anki) + VCs
There are exceptions such as DeepMind.
If you aim at a big exit, leave academia.

**Management:** There are 1001 things to do, so take a full-time CEO.

**Customers:** Selling your new technology provides satisfaction.
But: Many customers → many problems.
A single failure can lead to threats and endless discussions.
Automotive suppliers are particularly harsh, being pressured by OEMs.

**Engineers vs Salespeople**: Salespeople promise what Engineers cannot deliver.
You are in the middle…

**Lead promises are broken**. 50% of the times "we will buy / want to buy your product"
is  broken. Main reasons:

(i)  not a ``decider" (AT Biz jargon:  "nicht satisfaktionsfähig")
(ii) sudden change of management (happens all the time)

**VCs** :   - They can afford to the risk that only 1/10 of their companies succeed.
           - They do not want moderate exits and can block them.
           - They can force you to sell in some cases.
           - Good VCs help you.

**Business Angels**: Are normally the better solution when aiming at moderate exits.

**Universities**: They are usually the IP holders and want to get equity.
            Oxford: 20-30%, even though
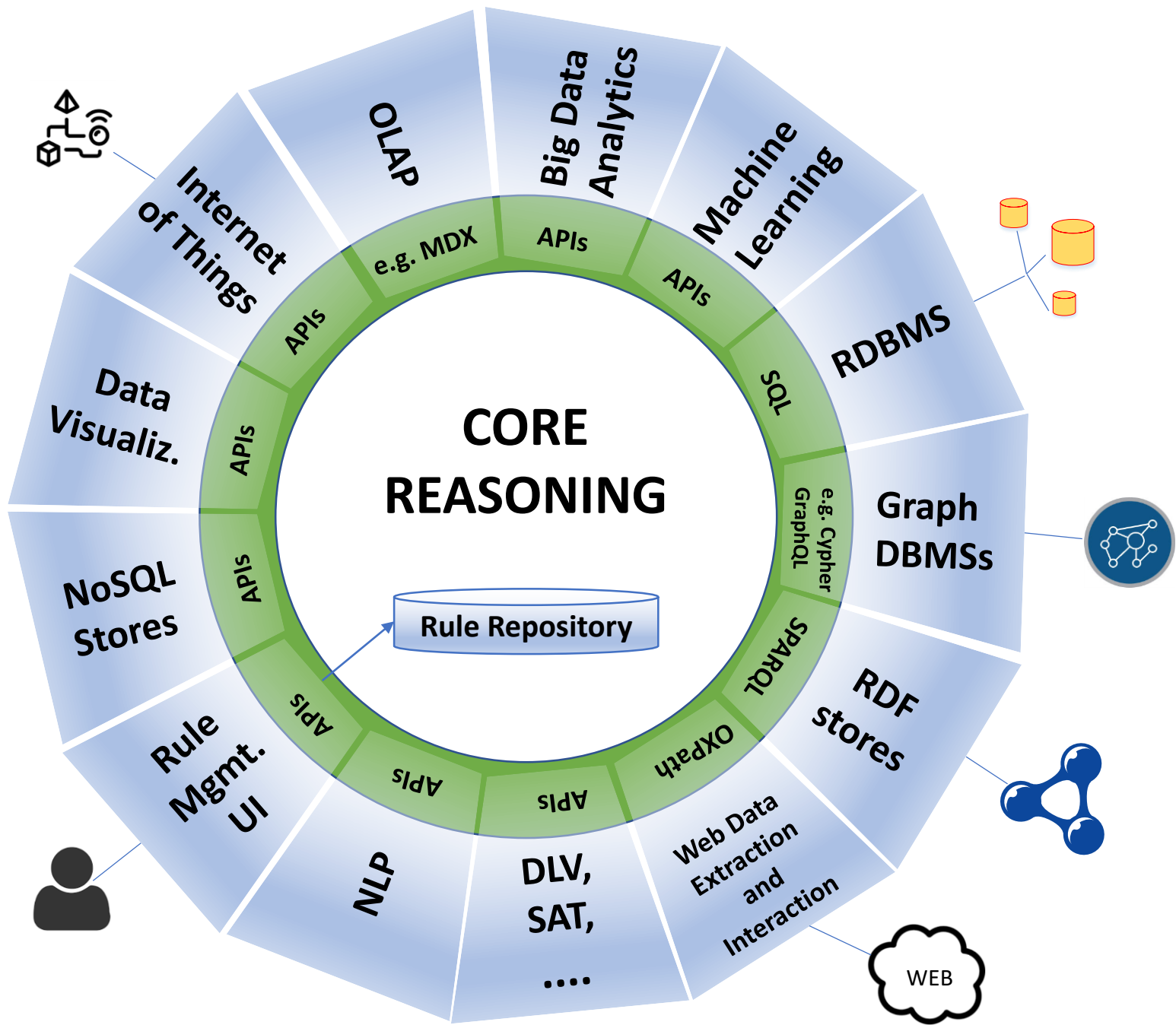            Austria??

**Advantages**:  Many very nice research problems arise from customer needs.

           Diverse activities let you "live longer".

           You prove that your research is useful. – better credibility,
            also  for publications.

           Additional income.

           Selling shares is tax efficient. In AT: 27.5% CGT vs 50% income tax.

# Enterprise Knowledge Graphs

**Facebook Knowledge Graph:** Social graph with people, places and things + information from Wikipedia

**Amazon Knowledge Graph:** Started as product categorization ontology

**Wolfram KB:** World facts + mathematics

**Factual:** Businesses & places

**Recruit Institute of Technology** (R.I.T): People, skills, recruiting

**Central Banks:** Company register – ownership graph

**Credit Rating Agencies …**

Thousands of medium to large size companies now want their own corporate knowledge graph. This not just for semantic indexing and search, but for advanced reasoning tasks on top of machine learning.

# Reasoning in Knowledge Graphs



EDB/ABox

Ontology / Rules

EDB+IDB

Many still think that DLs or graph databases suffice.  However:

Reasoning tasks are required that cannot be expressed by description logics, and cannot be reasonably managed by relational DBMS, nor by graph DBMS.

# Example: Wikidata Marriage Intervals

[Krötzsch  DL 2017]

Wikidata contains the statement :

**Taylor was married to Burton starting from 1964 and ending 1974**

This can be represented in relational DB or Datalog-notation by :

```
married(taylor,burton,1964,1974)
```

Symmetry rule for marriage intervals in Datalog:

$$\forall\, u,v,x,y.\ \texttt{married(u,v,x,y)} \rightarrow \texttt{married(v,u,x,y)}$$

**This cannot be expressed in DLs!**

Note: In what follows, we will often omit universal quantifiers.

# Example: My Creditworthiness

# Example: My Creditworthiness



up to £10,000

£8,500

£12,000

up to EUR 10,000

up to EUR 20,000

£500

£ 8,000

£ 12,500

EUR 14,000

# My Explanation

A machine-learning program has "reasonably" learned:

*People who live in a joint household with someone who does not pay their bills are likely to fail repaying their own debts.*

This ethically questionable rule was applied to
incomplete and wrong data:

• Before I bought the house there was a tenant who indeed did not pay his bills (*tons of unpaid bills & overdue notices in my mailbox*).

• The tenant had moved out before I moved in, but the Credit Rating Agency did not know, and simply assumed he still lived there.

ML should be complemented and, where necessary,  overruled by  domain-specific "expert rules"  that express domain knowledge:

- **A new house-owner is most likely unrelated to a previous tenant.**

- **If a house is bought by somebody who did not live there previously, and now lives there, then the previous occupiers have most likely moved out. [→Verify!]**

- **If someone has closed their bank account without opening a new one  then it is likely that the person has moved out of the country.**

  **...**

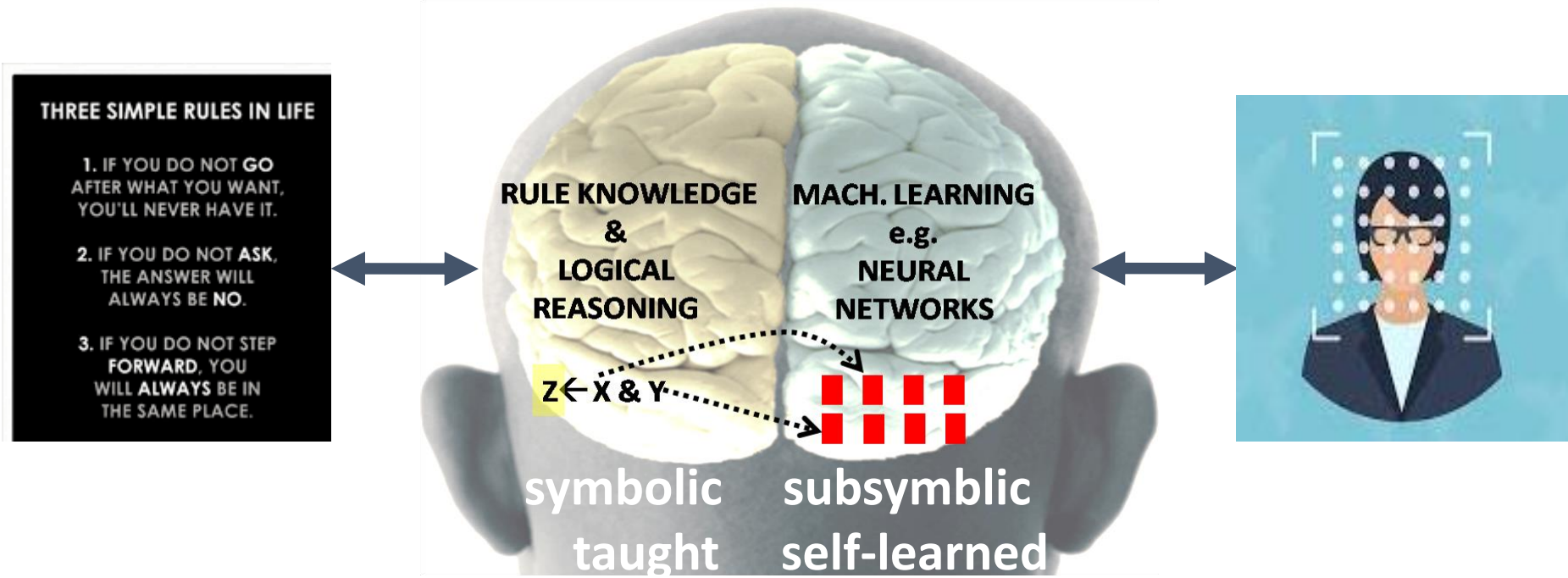Automatically accessing outside sources such as the Land Register and/or Social Networks may help.

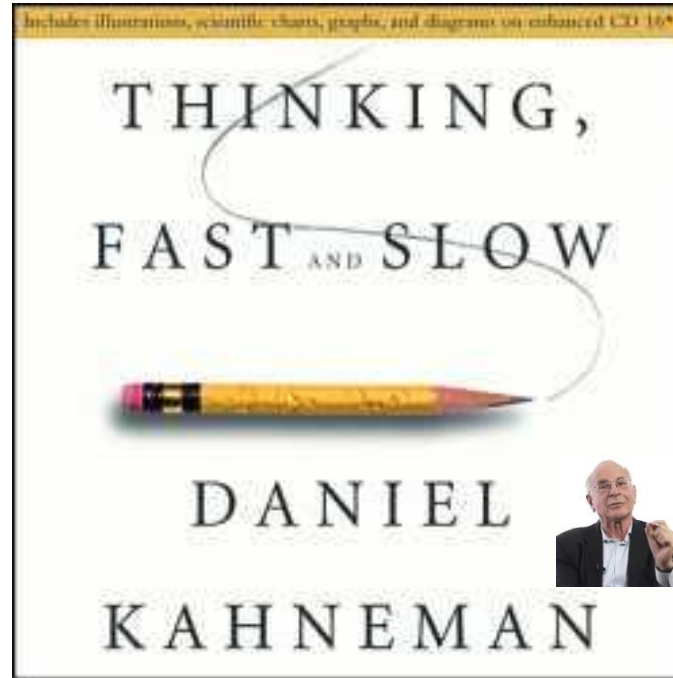→ Data extraction from external sources is a requirement for KGMS.

# Knowledge Graph Management Systems

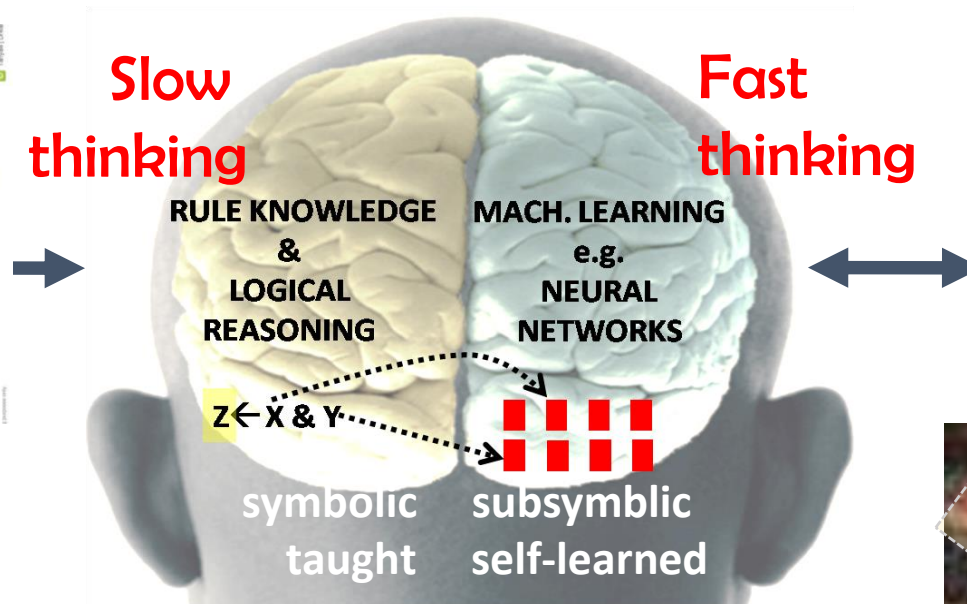KGMS combine the power of rule-based reasoning with machine learning over Big Data:

## KGMS = KBMS + Big Data + Analytics

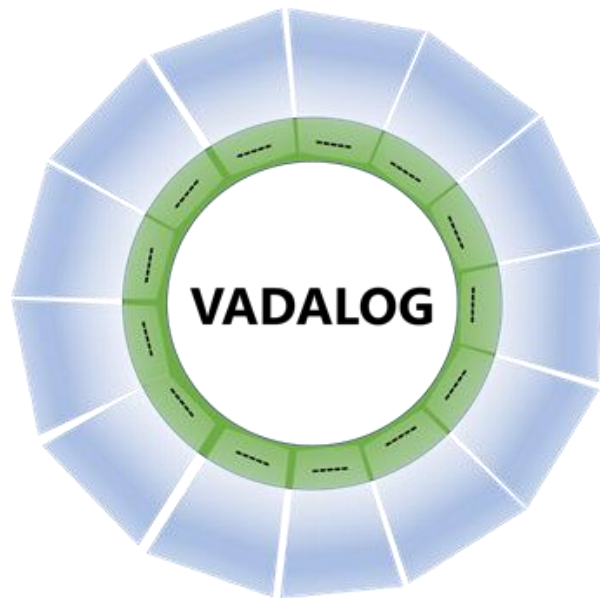**Misusing the lateralization thesis for illustration**

THINKING, FAST AND SLOW
DANIEL KAHNEMAN

Slow thinking

Fast thinking

RULE KNOWLEDGE & LOGICAL REASONING

MACH. LEARNING e.g. NEURAL NETWORKS

$Z \leftarrow X \ \& \ Y$

symbolic taught

subsymblic self-learned

Grandma: "Fly agarics are poisonous mushrooms. If you eat a poisonous mushroom, you may die".

Yikes, a fly agaric!

# Vadalog KGMS Being Built at Oxford



VADALOG



Yavor Nenov
Tim Furche
Stéphane Reissfelder
Vishal Chakhraborti
Emanuel Sallinger
Lianlong Wu
Luigi Bellomarini
Georg Gottlob
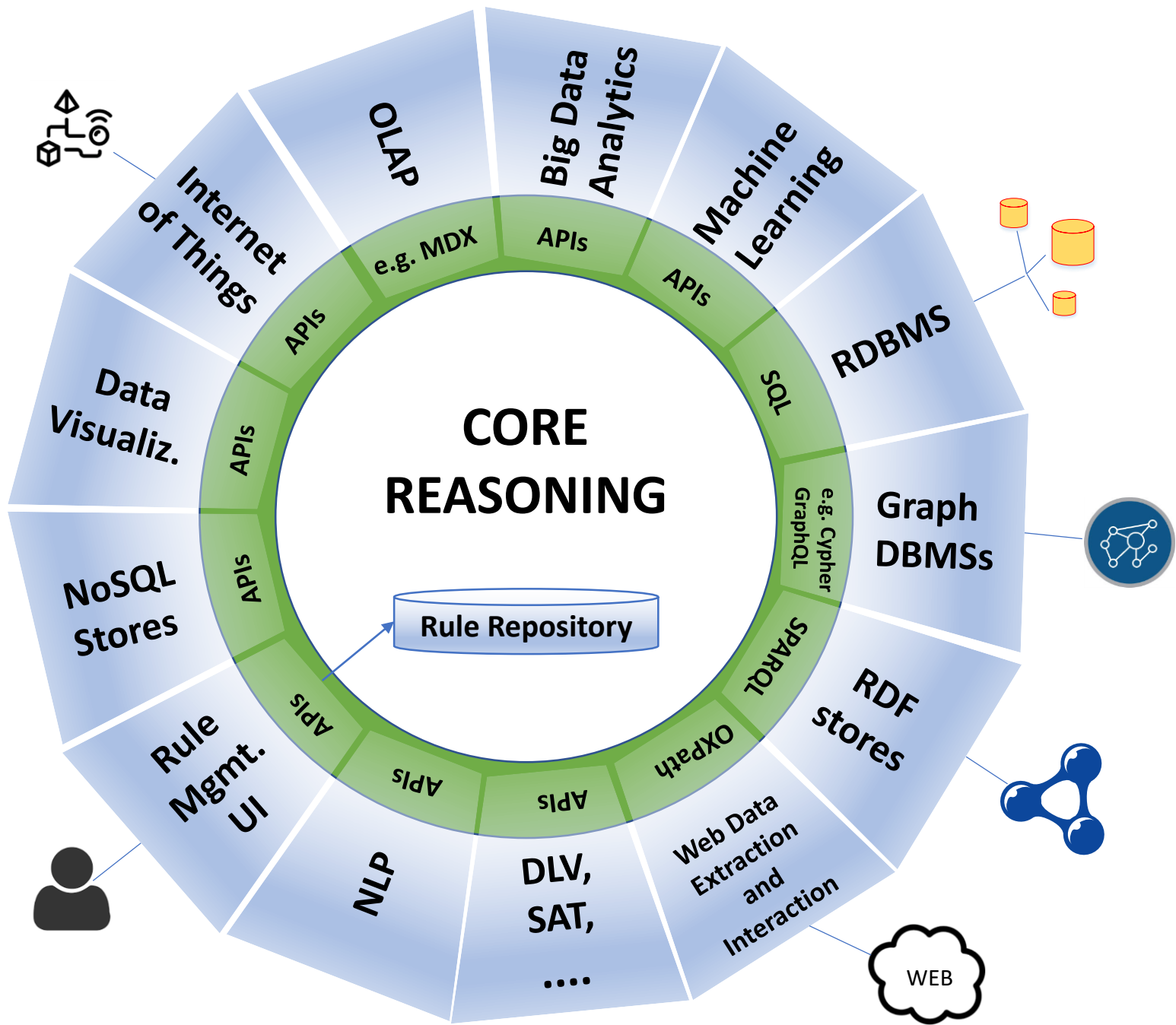Andrey Kravchenko
Julia Wiedmann
Ruslan Fayzrakhmanov
Evgeny Sherkhonov
Gerald Berger

Current Team Members

- VADA = **V**alue-**A**dded **DA**ta

- General architecture

- The core reasoning language "Warded Datalog" and its extensions
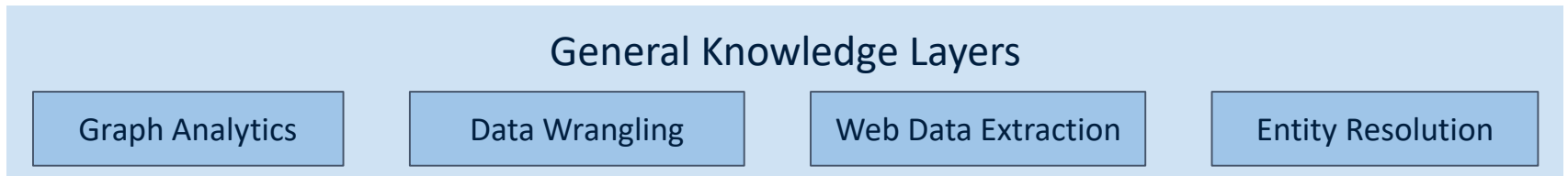
- Connectivity: Some plug-ins

# Knowledge Layers

**Vertical-specific Knowledge Layers**

| Logistics | Banking & Finance | Oil & Gas | Media Intel | Life Sciences | … | … |

**General Knowledge Layers**

| Graph Analytics | Data Wrangling | Web Data Extraction | Entity Resolution |

## Core Reasoning Engine
Strong performance and Expressiveness, Graph Navigation,
+    Integrations with Machine Learning & Enterprise Databases

DeepReason.ai