# The iMeMex Dataspace Management System: Architecture, Concepts, and Lessons Learned
## (Invited Tutorial)

Jens Dittrich

SAARLAND UNIVERSITY

COMPUTER SCIENCE

BNCOD 2009

# How it all started

- 2004:



Mac OS X



Linux



Windows

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# What is Personal Information?

Files&Folders

Calendar

Email plus Attached Files

Pictures & Videos

Music

Web-sites

RSS/ATOM Feeds

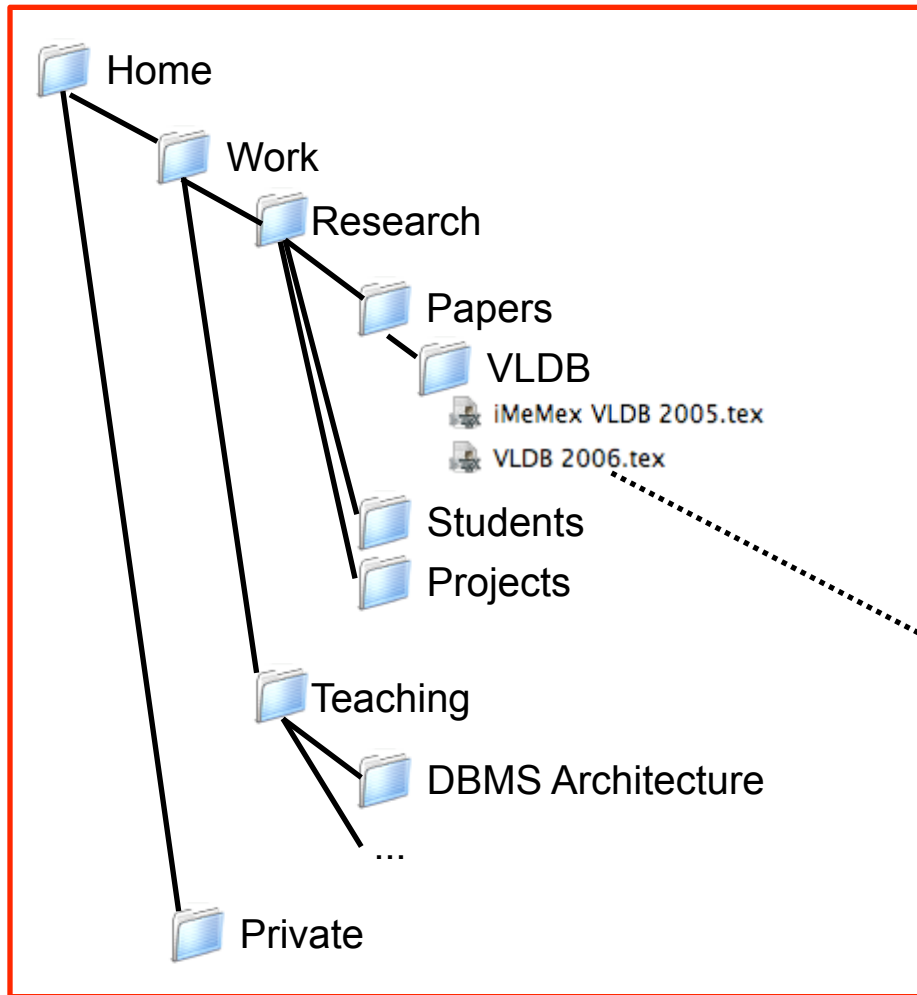# Problem 1: Users Store Stuff on Devices

- C: or network drive T:

- copy from C: to T:

- copy from C: to USB drive

- download pictures from digital camera to laptop

- download stuff from the Internet to laptop

- replicate data for backups between devices

- **Observation**: user knows about physical devices.

Users perform physical data management.

# Problem 2: Information Silos

# Problem 3: Artifical File Boundaries

Home
Work
Research
Papers
VLDB
iMeMex VLDB 2005.tex
VLDB 2006.tex
Students
Projects
Teaching
DBMS Architecture
...
Private

The outside world

- How to query all VLDB papers citing one of "Klaus Dittrich" papers from the late nineties?

- How to query all Teaching material citing "Klaus Dittrich" in any "architecture" lecture?

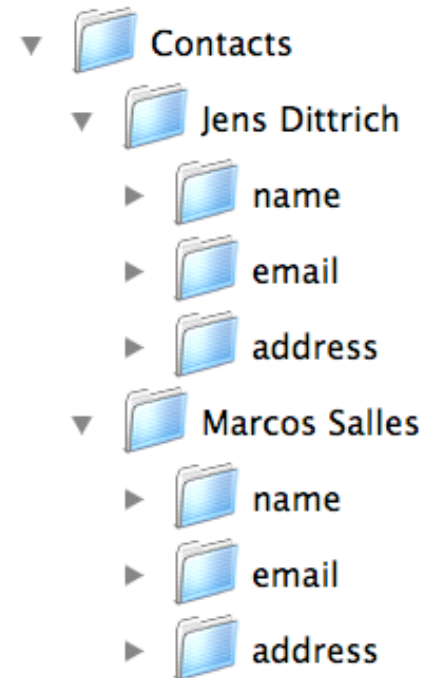- How to find all emails from those persons I cited in any paper I have published in 2005 or 2006?

```
\documentclass{vldb}
\title{iDM: A Unified ...}
\abstract{Personal Information...}
\begin{document}
\section{Introduction}
Personal Information...
...
\subsection{The Problem}
... basic concepts in Section~\ref{sec:preliminaries} ...
\section{Preliminaries}
\label{sec:preliminaries}
Intentional data can also...
\end{document}
```

The inside world

**Problem**: There is a gap between the outside and the inside structure.

# Problem 3: Artifical File Boundaries or: Data Format versus Data Model

```xml
<contacts>
  <contact id=1>
    <name>Jens Dittrich</name>
    <email>jens.dittrich@cs.uni-sb...</email>
    <address>ETH...</address>
  </contact>
  <contact id=2>
    <name>Marcos Salles</name>
    <email>marcos.salles@cornell...</email>
    <address>ETH...</address>
  </contact>
    ....


</contacts>
```
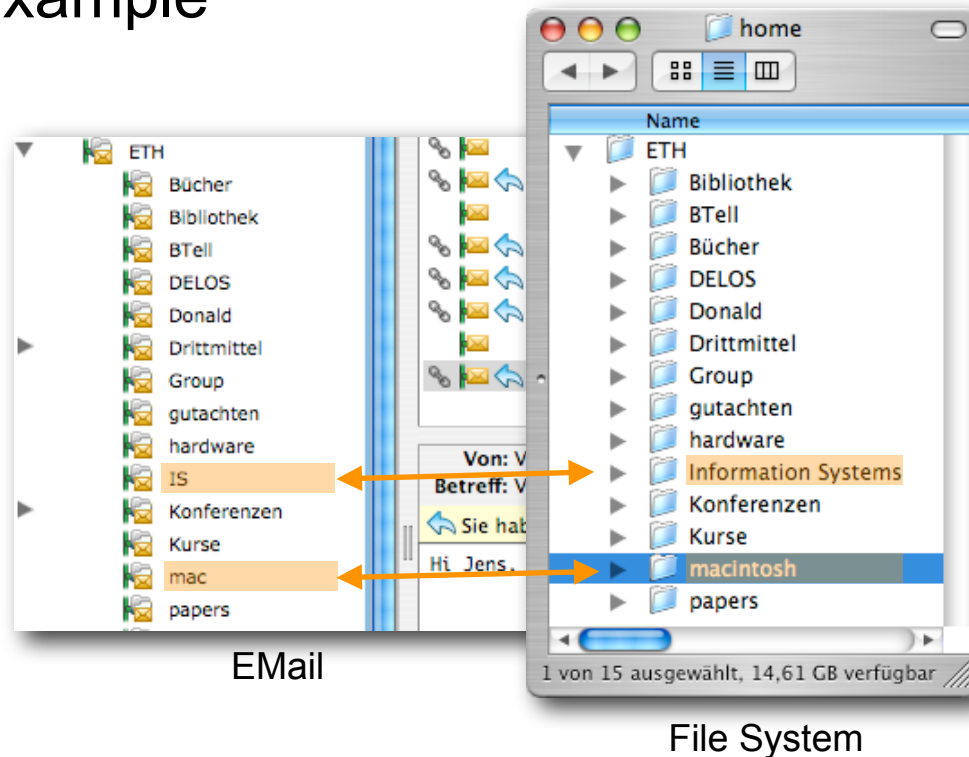


- Where does the folder-hierarchy end and the XML start?
- It is up to the user to define the boundary.

**Problem**: Same data model but different formats/representations.

# Problem 4: Repeated Folder Hierarchies

- ## Example



EMail

File System

- Similar hierarchies in multiple places
  - local desktop disk
  - local laptop disk
  - network drive
  - email folders
  - bookmarks

This is a mix physical data management and manual schema mapping.

SAARLAND UNIVERSITY
COMPUTER SCIENCE

# Problem 5: Finding

- How to find stuff in this device and format mess?
- just keywords?
- how to query the hidden structures?
- how to link similar structures, e.g. hierarchies?
- searching versus querying

# PIM Hell

**Today**: Users have to perform too many physical data managing and schema mapping tasks.

1. Users store stuff on devices, e.g. PC, Laptop, iPod, cell phone, USB stick, server; copy among devices, etc.

2. application silos, structural content hidden inside files

3. artificial file boundaries: inside versus outside hierarchies

4. similar folder hierachies on different devices
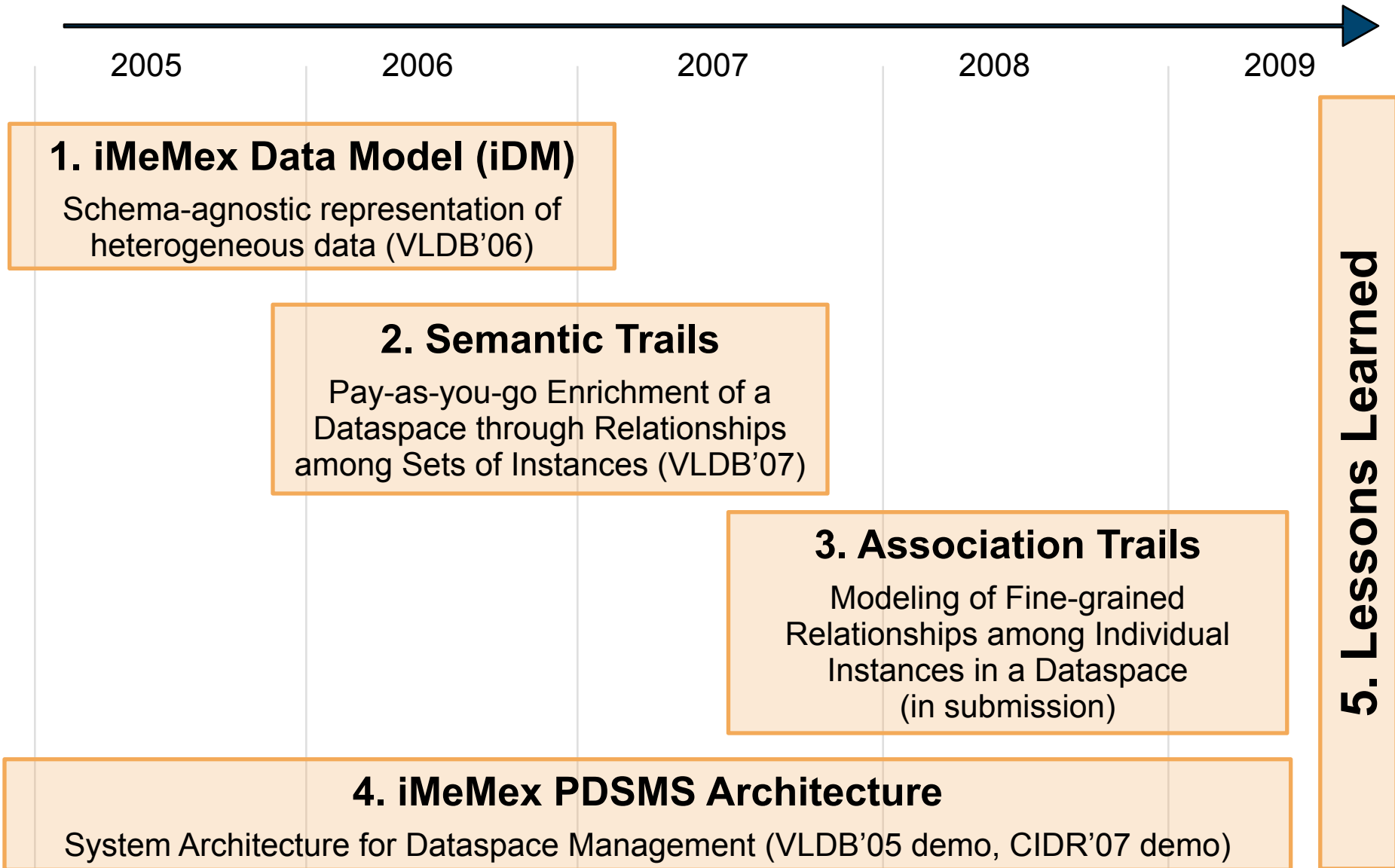
5. finding

6. etc.

# PIM Heaven

> **Tomorrow**: Users should only do logical data management and do not worry (too much) about schemas.
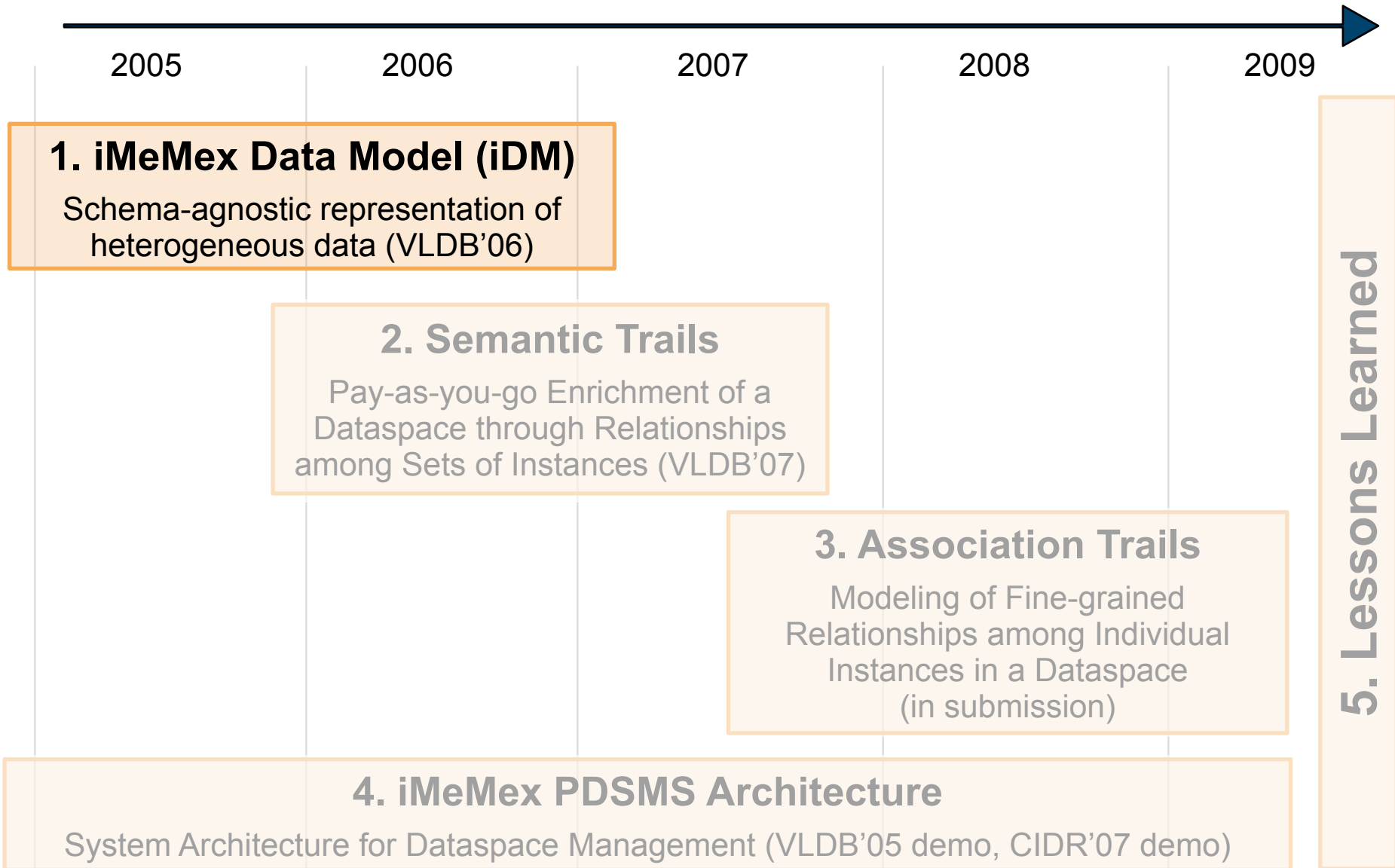
- **Goals**
  - get rid of physical data management
  - logical granularity should be independent from the physical unit
  - **integrate** data without all the hazzle of complex schema integration
  - allow for powerful search **and** query facilities
  - ...
- **Challenge**: build a system that is able to do that...
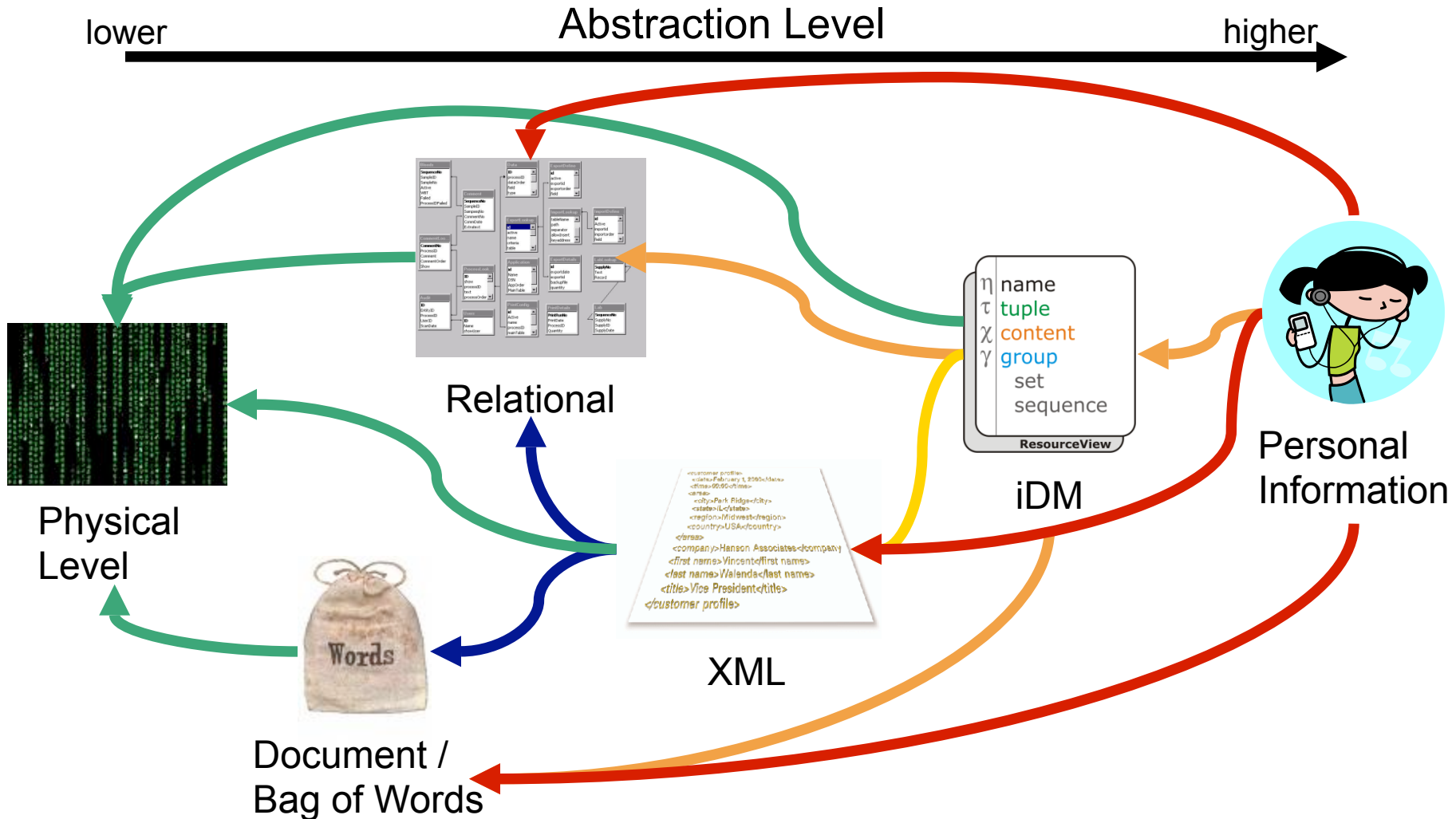
# Outline

2005      2006      2007      2008      2009

**1. iMeMex Data Model (iDM)**

Schema-agnostic representation of heterogeneous data (VLDB'06)

**2. Semantic Trails**

Pay-as-you-go Enrichment of a Dataspace through Relationships among Sets of Instances (VLDB'07)

**3. Association Trails**

Modeling of Fine-grained Relationships among Individual Instances in a Dataspace (in submission)

**4. iMeMex PDSMS Architecture**

System Architecture for Dataspace Management (VLDB'05 demo, CIDR'07 demo)

**5. Lessons Learned**

SAARLAND
UNIVERSITY

COMPUTER SCIENCE

# Outline

2005      2006      2007      2008      2009

## 1. iMeMex Data Model (iDM)

Schema-agnostic representation of heterogeneous data (VLDB'06)

## 2. Semantic Trails

Pay-as-you-go Enrichment of a Dataspace through Relationships among Sets of Instances (VLDB'07)

## 3. Association Trails

Modeling of Fine-grained Relationships among Individual Instances in a Dataspace (in submission)

## 4. iMeMex PDSMS Architecture

System Architecture for Dataspace Management (VLDB'05 demo, CIDR'07 demo)

## 5. Lessons Learned

SAARLAND UNIVERSITY COMPUTER SCIENCE

# Personal Information is...

- Non-schematic, heterogeneous collections with no formal schema

- Serialized in hundreds of file formats and encodings

- Organized in arbitrary graphs (inside and outside of files)

- Distributed among different data sources

- Potentially infinite (e.g. RSS, ATOM, email streams)
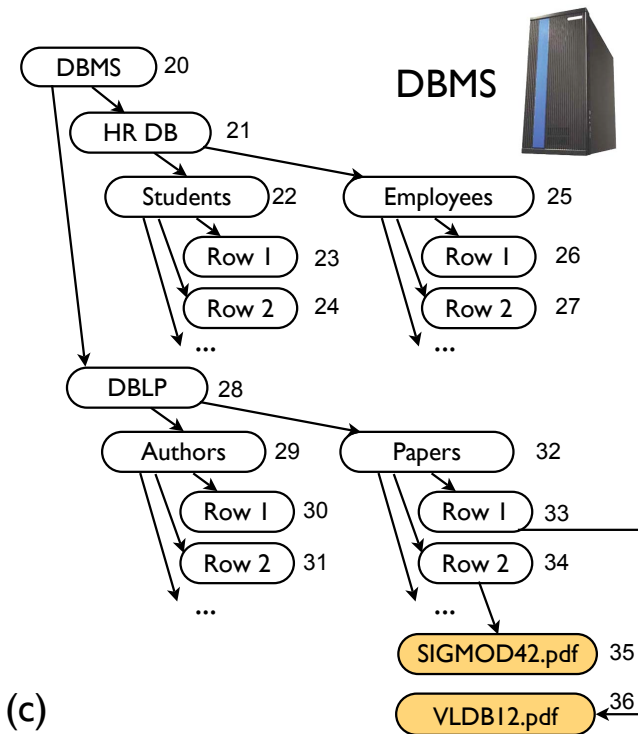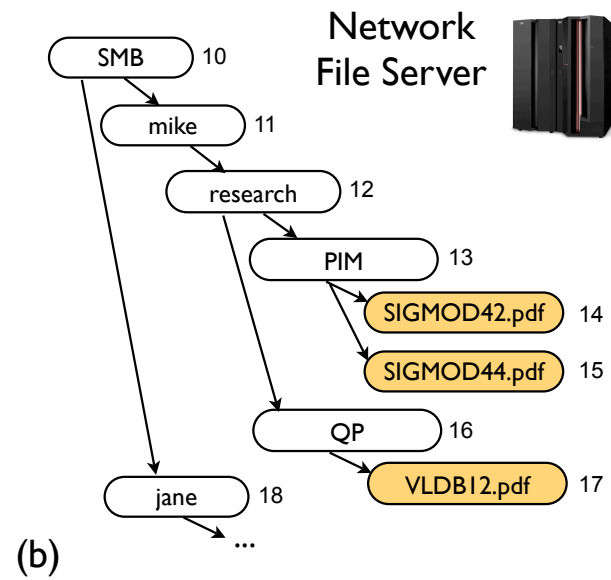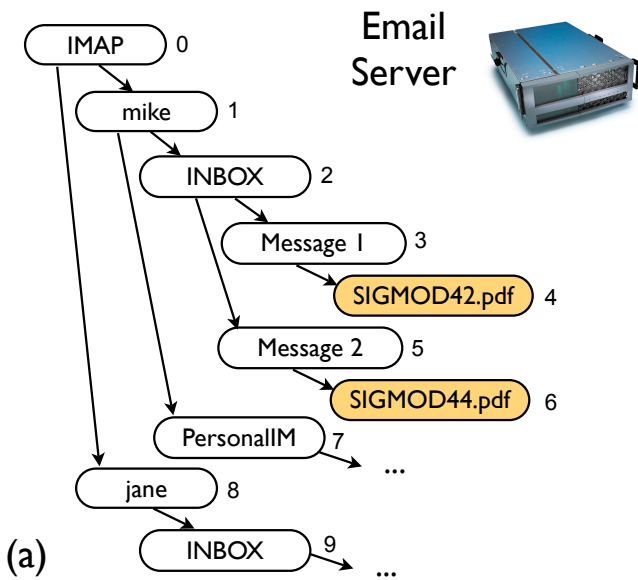
# Data Models for Personal Information

# iDM: Representing Information in Personal Dataspaces

- **Our approach:** represent all personal information into a common data model and offer a unified query-and-search service

- **Applications:**

  - A powered-up shell → paths and keywords across filesystems, email, databases, outside and inside of files, e.g. `//projects/main.tex/section/subsection["mike"]`

  - A powered-up search application → return not only files as results, but elements at arbitrary granularity, e.g. bibliographic references
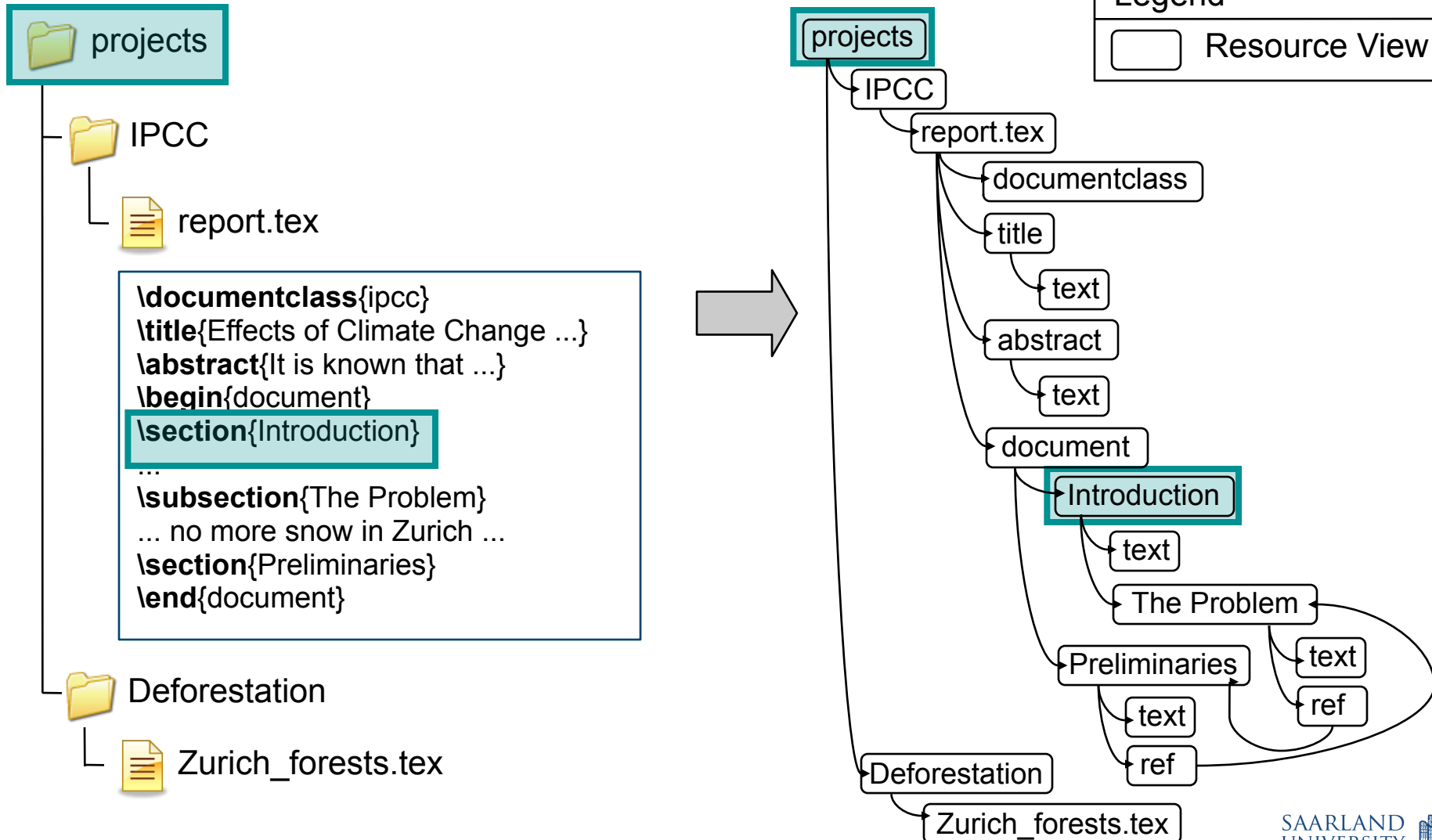
# iDM Core Idea: Lazily Computed Graph

- Nodes and Edges are lazily computed
- Each node is termed *Resource View*

**(a) Email Server**

- IMAP — 0
- mike — 1
- INBOX — 2
- Message 1 — 3
- SIGMOD42.pdf — 4
- Message 2 — 5
- SIGMOD44.pdf — 6
- PersonalIM — 7 ...
- jane — 8
- INBOX — 9 ...

**(b) Network File Server**

- SMB — 10
- mike — 11
- research — 12
- PIM — 13
- SIGMOD42.pdf — 14
- SIGMOD44.pdf — 15
- QP — 16
- VLDB12.pdf — 17
- jane — 18 ...

**(c) DBMS**

- DBMS — 20
- HR DB — 21
- Students — 22
- Row 1 — 23
- Row 2 — 24 ...
- Employees — 25
- Row 1 — 26
- Row 2 — 27 ...
- DBLP — 28
- Authors — 29
- Row 1 — 30
- Row 2 — 31 ...
- Papers — 32
- Row 1 — 33
- Row 2 — 34 ...
- SIGMOD42.pdf — 35
- VLDB12.pdf — 36

**(d) Laptop**

- home — 40
- mike — 41
- papers — 42
- PIM — 43
- SIGMOD42.pdf — 44
- SIGMOD44.pdf — 45
- QP — 46
- VLDB12.pdf — 47
- VLDB10.pdf — 48
- projects — 49
- PIM — 50
- SIGMOD42.pdf — 51

# iDM Removes the Inside-Outside File Boundary



```
\documentclass{ipcc}
\title{Effects of Climate Change ...}
\abstract{It is known that ...}
\begin{document}
\section{Introduction}
...
\subsection{The Problem}
... no more snow in Zurich ...
\section{Preliminaries}
\end{document}
```

**Legend**

Resource View

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# iDM Features: Lazy Computation

- Important: iDM is not a static model.

- Every component of every Resource View may be created on demand.

- Furthermore, every Resource View may be created on demand.

- This achieved by modeling a Resource view as a set of get*-methods:

```
Interface ResourceView {
        getNameComponent(): return η
        getTupleComponent(): return τ
        getContentComponent(): return χ
        getGroupComponent() : return γ
}
```
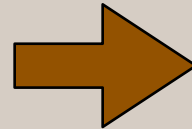
**Important**: It is up to the dataspace system to decide when the result to a get*-method is materialized.

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# iDM Features: Lazy Computation Examples

```
Interface ResourceView {
    getNameComponent(): return η
    getTupleComponent(): return τ
    getContentComponent(): return χ
    getGroupComponent() : return γ
}
```

- getContent
    - system retrieves web page from a remote server
    - or: system dynamically generates a html page
    - or: system returns an already cached web page
    - etc.

- getGroup
    - system calls getContent, extracts structural information, returns it as an iDM subgraph
    - or: system processes a query and returns result as iDM subgraph
    - or: system calls a web service and returns result as iDM subgraph
    - or: system returns an already cached group component
    - or: system retrieves group component from a remote server

**Important**: the dataspace system has to make decisions on resource view materialization.

UNIVERSITY
COMPUTER SCIENCE

# iDM Features: Use-case Active XML

**Active XML**

Proposed by Abiteboul et.al. PODS 04, SIGMOD 04, PODS 05, etc.

```
<dep>
  <sc>web.server.com/GetDepartments()</sc>
</dep>
```

```
<dep>
    <sc>web.server.com/GetDepartments()</sc>
    <deplist>
        <entry>
            <name>Accounting</name>
        </entry>
        ...
    </deplist>
</dep>
```

(1) Original XML document

(2) Same XML document after calling web service

**iDM**

How to use iDM to achieve the same effect:

$$\gamma_i^{AXML} = \left( \varnothing, \langle V_j^{sc}[, V_k^{scresult}] \rangle \right)$$

SAARLAND UNIVERSITY
COMPUTER SCIENCE

# iDM Features: Built-in Stream Support



- **Infinite components may occur in 3 places of a resource view**
  - (1) content component (stream of characters)
    - Example: video and audio stream broadcast over the network
  - (2) set or (3) sequence of the group component (stream of Resource Views)
    - Examples
      - any data stream
      - pub/sub system
      - sensor data

# iDM Use-case: Email

- Consider all emails routed to address jens.dittrich at cs..
  Two options to model this using iDM:

  1. Option: Model the state:

  - $$\gamma_i^{\text{INBOX State}} = (\{\}, \langle V_{\mathbf{q}_1}^{\text{message}}, \ldots, V_{\mathbf{q}_n}^{\text{message}} \rangle)$$

  - Note: the INBOX represents a window query = some state is preserved.

  - The state of that query is equal to the list of messages contained in the INBOX (shedding is performed by user or spam-filter).

  - Messages may be retrieved multiple times.

  2. Option: Model the stream:

  - $$\gamma_i^{\text{INBOX message stream}} = (\{\}, \langle V_{\mathbf{q}'_1}^{\text{message}}, \ldots, V_{\mathbf{q}'_n}^{\text{message}} \rangle_{n \to \infty})$$

  - Stateless approach
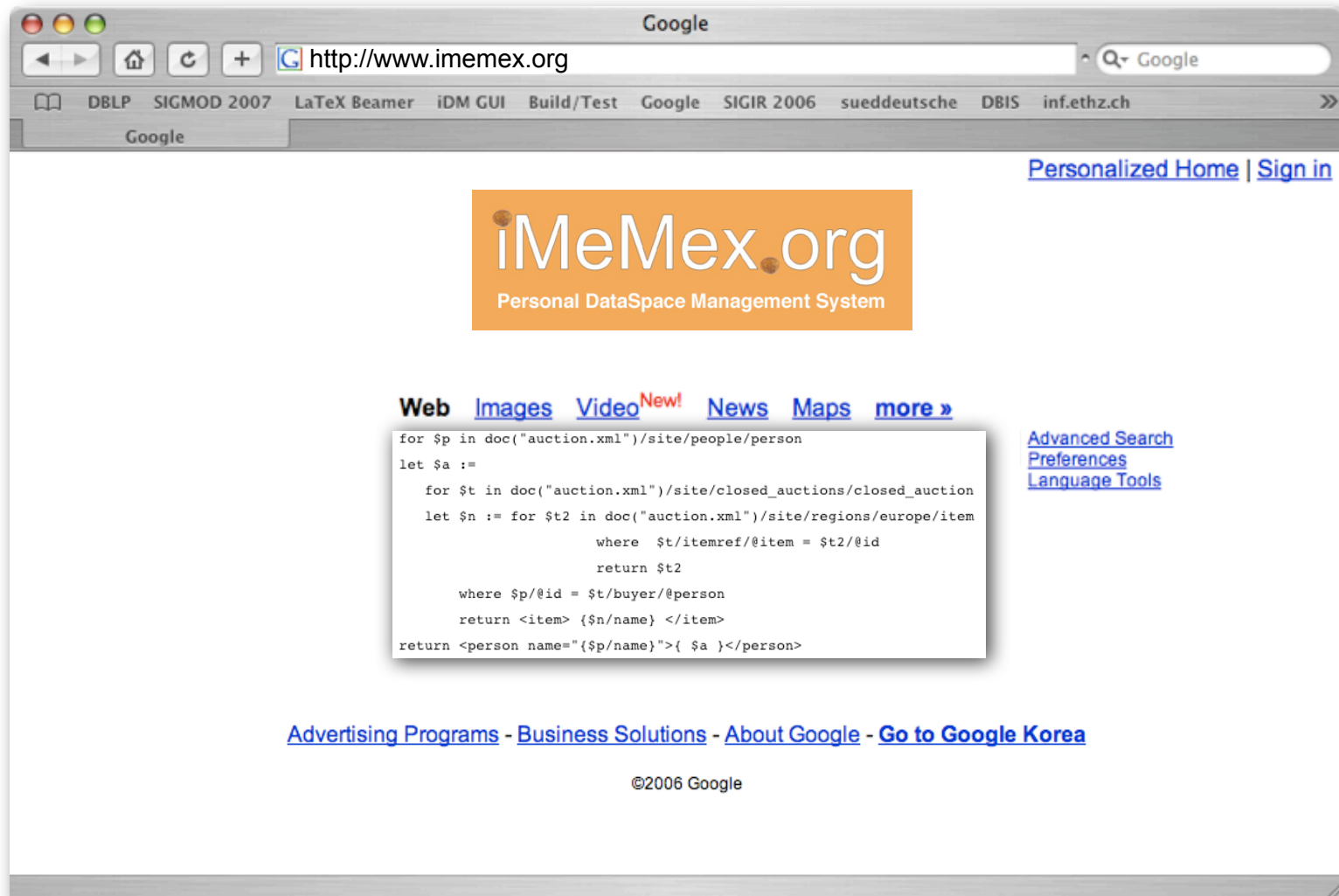
  - Messages cannot be retrieved a second time.

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# Personal Information Features vs. Data Models

| Support for Personal Data | | Bag of Words | Relational | XML | iDM |
|---|---|---|---|---|---|
| | Non-schematic data | ✅ | 🚫 | ✅ | ✅ |
| | Serialization agnostic | ✅ | ✅ | 🚫 | ✅ |
| | Support for Graph data | 🚫 | Specific schema | Extension: XLink/ XPointer | ✅ |
| | Support for Lazy Computation | 🚫 | View mechanism | Extension: ActiveXML | ✅ |
| | Support for Infinite data | Extension: Document streams | Extension: Relational streams | Extension: XML streams | ✅ |

SAARLAND UNIVERSITY
COMPUTER SCIENCE

# How to Query the iDM Dataspace? Like this?

# Or like this?

# iQL: Towards a Dataspace Query Language

- Language Requirements
  - simple and expressive at the same time
  - centered around keyword search
  - should have structural constraints
  - algebraic operations (joins)
  - support updates and inserts.

- Existing search&query languages
  - keyword search: no structural constraints, too leightweight
  - SQL: too complex, too much focussed on relational model
  - XPath : good on structural constraints, bad on keywords
  - XQuery: far too heavy

# Our Approach: iQL

- `Donald Knuth`
  returns all resource views containing both keywords "Donald" and "Knuth"

- `"Donald Knuth"`
  returns all resource views containing the phrase "Donald Knuth"

- `[size > 42000 and lastmodified < yesterday()]`
  returns those resource views having a tuple component attribute greater than 42000 and a lastmodified date before yesterday.

- `//PIM//Introduction[class="latex_section"]`
  returns every resource view named "Introduction" of class "latex_section" that is indirectly related to a resource view named "PIM".

- `//OLAP//[class="figure" and "Indexing time"]`
  first, selects resource views that are indirectly related to a resource view named "OLAP". In addition, all results have to be of resource view class "figure" and have to contain the phrase "Indexing time".

- In the IR-community a related approach was proposed restricted to XML retrieval: NEXI (Narrowed Extended XPath), Trotman and Sigurbjörnsson, INEX 2004

- However, NEXI is simply not powerful enough.

SAARLAND
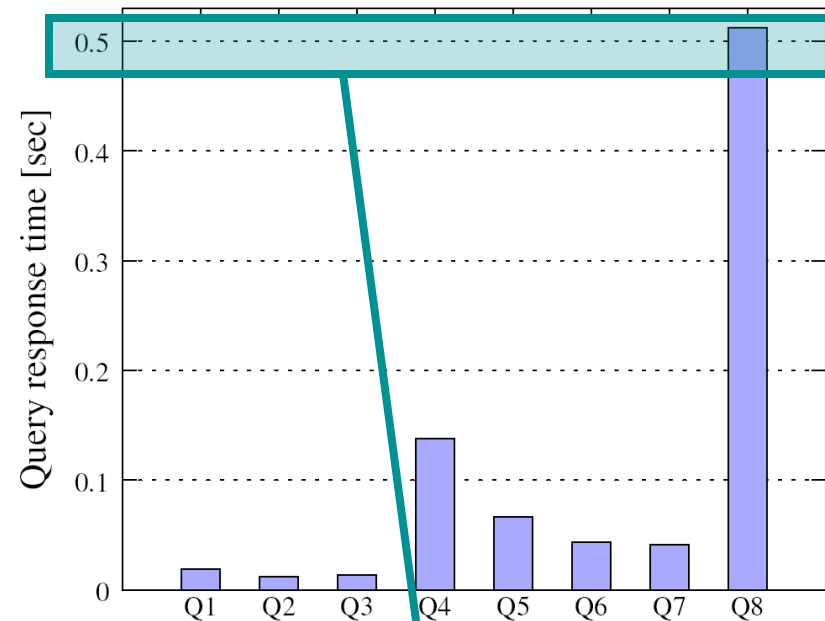UNIVERSITY
COMPUTER SCIENCE

# iQL: iMeMex Query Language

- **Core idea:** intuitive keyword&path language to search iDM graphs

- **Examples**

  - `global warming zurich`

  - `celsius > 10 and region = "ZH"`

  - `//inbox/IPCC//*.pdf`

  - `//Temperatures/*[region = "ZH"]`
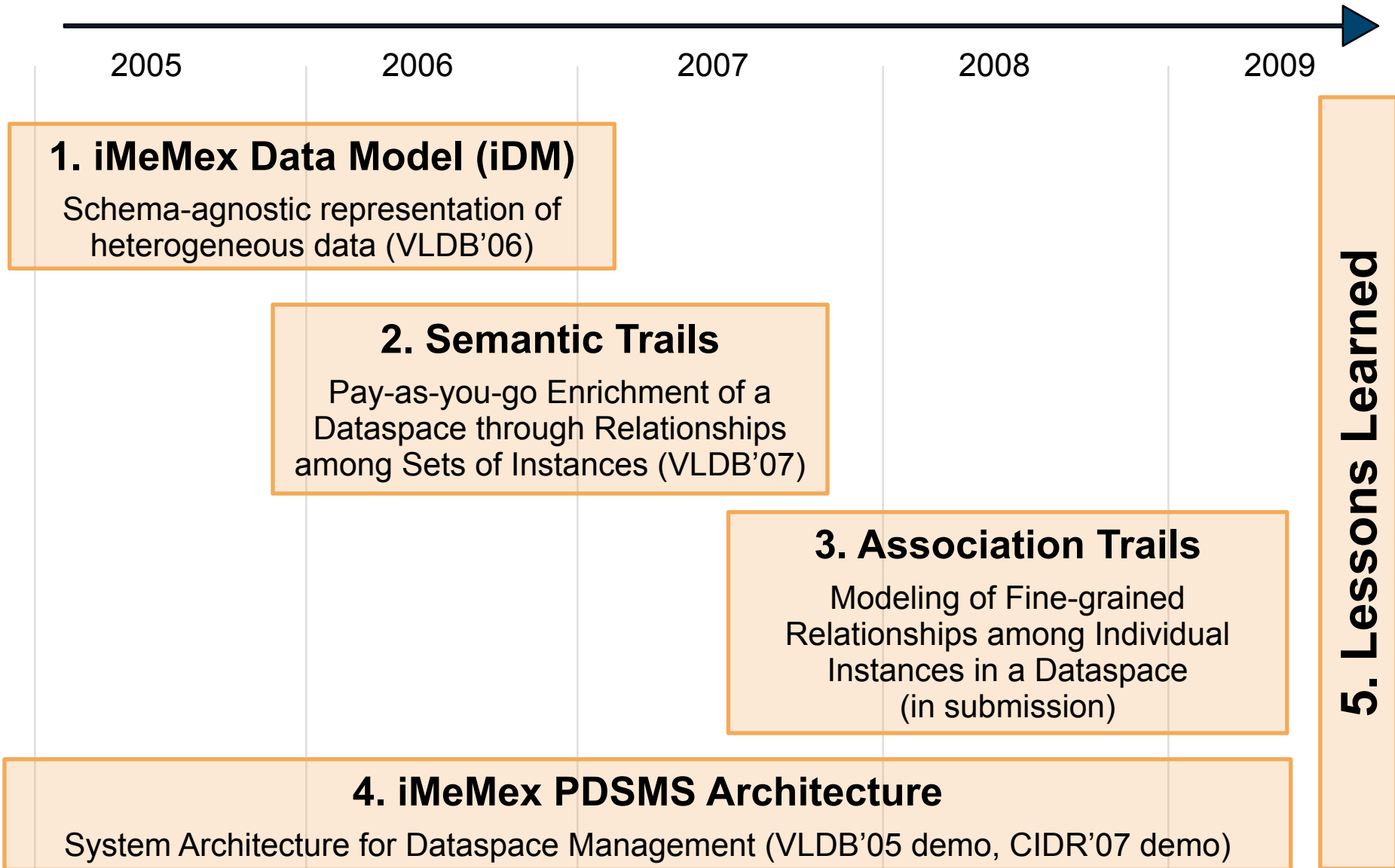
# Evaluation of iDM and iQL in iMeMex

- Personal dataset from filesystem and email
- Indexing of iDM graphs with inverted lists & group replica

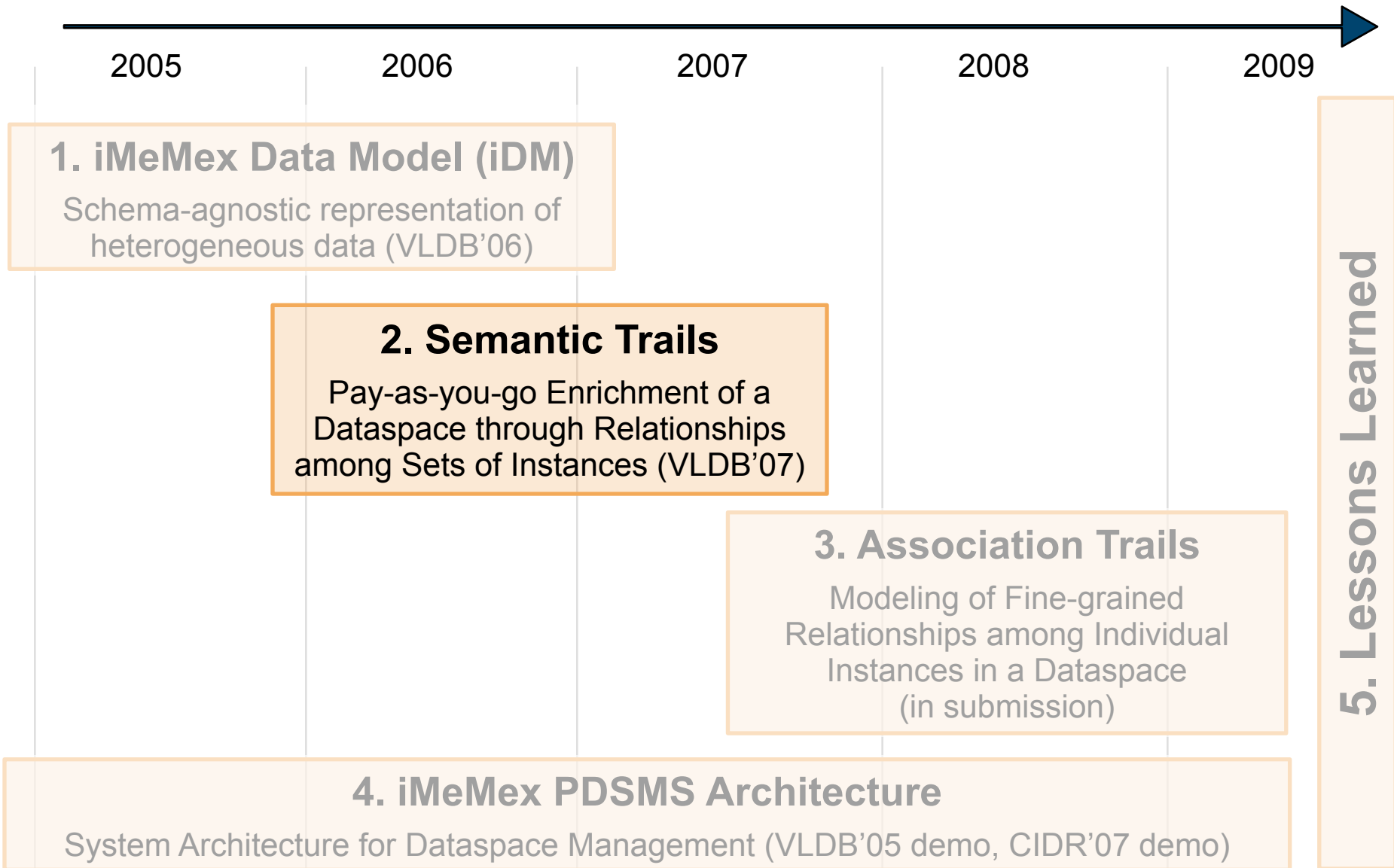| | iQL Query expression | # of Results |
|---|---|---|
| Q1 | `database` | 941 |
| Q2 | `database tuning` | 39 |
| Q3 | `size > 420000 and lastmodified < 12.06.2005` | 88 |
| Q4 | `//papers//*Vision/*["Franklin"]` | 2 |
| Q5 | `//VLDB200?//?onclusion*/*["systems"]` | 2 |
| Q6 | `//VLDB2005//*["documents"] ∪`<br>`        //VLDB2006//*["documents"]` | 31 |
| Q7 | `join( //VLDB2006//*[class="texref"] as A,`<br>`//VLDB2006//*[class="environment"]//figure* as B,`<br>`A.name=B.tuple.label )` | 21 |
| Q8 | `join( //*[class = "emailmessage"]//*.tex as A,`<br>`//papers//*.tex as B, A.name = B.name )` | 16 |

Interactive query times, but variance resulting from implicit joins in path expressions

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# Outline



2005    2006    2007    2008    2009

**1. iMeMex Data Model (iDM)**

Schema-agnostic representation of heterogeneous data (VLDB'06)

**2. Semantic Trails**

Pay-as-you-go Enrichment of a Dataspace through Relationships among Sets of Instances (VLDB'07)

**3. Association Trails**

Modeling of Fine-grained Relationships among Individual Instances in a Dataspace (in submission)

**4. iMeMex PDSMS Architecture**

System Architecture for Dataspace Management (VLDB'05 demo, CIDR'07 demo)

**5. Lessons Learned**

SAARLAND UNIVERSITY

COMPUTER SCIENCE

# Outline

2005     2006     2007     2008     2009

**1. iMeMex Data Model (iDM)**

Schema-agnostic representation of heterogeneous data (VLDB'06)

**2. Semantic Trails**

Pay-as-you-go Enrichment of a Dataspace through Relationships among Sets of Instances (VLDB'07)

**3. Association Trails**

Modeling of Fine-grained Relationships among Individual Instances in a Dataspace (in submission)

**4. iMeMex PDSMS Architecture**

System Architecture for Dataspace Management (VLDB'05 demo, CIDR'07 demo)

**5. Lessons Learned**

SAARLAND UNIVERSITY
COMPUTER SCIENCE

# Problem: Lack of Unified View of Personal Data Sources

Query

What is the impact of global warming in Zurich?

? ? ? ?



Systems

Data Sources

Laptop

Email Server

Web Server

DB Server

# Solution 1: Use a Search Engine

**Query**

`global warming zurich`

**Search Engine**



**Data Sources**

text, links — Laptop

text, links — Email Server

text, links — Web Server

text, links — DB Server

**Drawback: Query semantics are not precise!**

UNIVERSITY
COMPUTER SCIENCE

# Solution 2: Use an Information Integration System



**Temperature, CO₂, and Sunspots**

`//Temperatures/*[region = "ZH"]`  Query

**Information**

**Drawback: Too much effort to provide schema mappings!**

GLAV [AAAI00], P2P (e.g. [SIGMOD04])

missing schema mapping | missing schema mapping | schema mapping | schema mapping

Laptop | Email Server | Web Server | DB Server

Data Sources

SAARLAND UNIVERSITY
COMPUTER SCIENCE

# Research Challenge:
# Is There an Integration Solution in-between These Two Extremes?

global warming zurich

**Temperature, CO₂, and Sunspots**



**Personal Dataspace Management System**

**Pay-as-you-go Information Integration**

text, links

full-blown schema mappings

**Data Sources**

text, links

text, links

text, links

text, links

**Data Sources**

Laptop

Email Server

Web Server

DB Server

Dataspace Vision by Franklin, Halevy, and Maier [SIGMOD Record 05]

SAARLAND UNIVERSITY
COMPUTER SCIENCE

# Schema-first vs. Dataspaces (From Mike Franklin's talk)

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# iTrails: Pay-as-you-go Definition of Relationships among Sets of Instances

- **Step 1:** Provide a search service over all the data
  - Use a general graph data model → iDM subset
  - Works for unstructured documents, XML, and relations

- **Step 2:** Add integration semantics via hints (semantic trails) on top of the graph
  - Works **across** data sources, not only between sources

- **Step 3:** If more semantics needed, go back to step 2

- **Impact:**
  - Smooth transition between search and data integration
  - Semantics added incrementally improve precision / recall

# iTrails: Definition of a Semantic Trail

- **Basic Form of a Semantic Trail**

Queries: iQL keyword and path expressions

$$Q_L [.C_L] \longrightarrow Q_R [.C_R]$$

Attribute projections

- **Intuition:** When I query for $Q_L [.C_L]$, you should also query for $Q_R [.C_R]$

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# Semantic Trails: Deep Web Bookmarks

`train home`



ZVV Reiseplaner — ZVV Richtig verkehrt.

Timetable Switzerland
+ door to door within canton Zurich (ZH)

| | | |
|---|---|---|
| From: | Station/Stop | eth uni |
| To: | Station/Stop | seilbahn rigiblick |
| Via(1): | Station/Stop | |
| Date: | Sa, 15.09.07 | ◄ ► Calendar |
| Time: | 19:04 | ⦿ Departure ○ Arrival |

Search connection | New query | More

- **Trail for a Bookmark:**
  "When **I** query for `train home`, you should also query the `TrainCompany`'s website"

```
train home →
    //trainCompany.com//*[origin="ETH Uni"
          and dest ="Seilbahn Rigiblick"]
```

Detailed view

| Station/Stop | Date | Time | Platform | Products | Comments |
|---|---|---|---|---|---|
| Zürich, ETH/Universitätsspital [i] [M] | 15.09.07 | dep 19:05 | | [🚋] Trm 9 | Trm Direction: Zürich, Hirzenbach |
| Zürich, Seilbahn Rigiblick [i] [M] | | arr 19:08 | | | |

Duration: 0:03; runs Sa
Hint: Departure/Arrival replaced by an equivalent station
[T] Tariff level*: 9; Zones*: 10; Short distance

SAARLAND UNIVERSITY COMPUTER SCIENCE

# Semantic Trails: Schema Equivalences

Employee

| empId | empName | salary |
|-------|---------|--------|

Person

| SSN | name | age | income |
|-----|------|-----|--------|

- **Trail for schema match on names:** "When I query for `Employee.empName`, you should also query for `Person.name`"

  ```
  //Employee//*.tuple.empName →
                  //Person//*.tuple.name
  ```

- **Trail for schema match on salaries:** "When I query for `Employee.salary`, you should also query for `Person.income`"

  ```
  //Employee//*.tuple.salary →
                  //Person//*.tuple.income
  ```

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# Semantic Trails: Multiple Hierarchies

global warming | zurich

```
global warming →
    //projects/IPCC//*
```

```
//projects/IPCC →
    //email/ClimateChange
```

```
zurich →
    region = "ZH"
```

projects

?

IPCC

report.tex

**\documentclass**{ipcc}
**\title**{Effects of Climate Change ...}
**\abstract**{It is known that ...}
**\begin**{document}
**\section**{Introduction}
...
**\subsection**{The Problem}
... no more snow in Zurich ...
**\section**{Preliminaries} ...
**\end**{document}

email

ClimateChange

Latest Temperature Data

Temperatures.dat

| date | city | region | celsius |
|------|------|--------|---------|
| 24-Sep | Bern | BE | 20 |
| 24-Sep | Uster | ZH | 15 |
| 25-Sep | Kloten | ZH | 14 |

Deforestation

Zurich_forests.tex

# Rewriting Queries with Semantic Trails: Multiple Match Coloring Algorithm



T₁: global warming → //projects/IPCC//*
T₂: //projects/IPCC → //email/ClimateChange
T₃: zurich → region = "ZH"

SAARLAND UNIVERSITY
COMPUTER SCIENCE

# How are Trails Created?

- Given by the user
  - Explicitly:
    - not unlikely for structural extensions
  - Via Relevance Feedback:
    - ask the user

- (Semi-)Automatically
  - Information extraction techniques
  - Automatic schema matching
  - Ontologies and thesauri (e.g., wordnet)
  - User communities (e.g., trails on gene data, bookmarks)

# (Semi-)Automatic Creation vs. Semantic Trail Rewrites

- Rewriting queries to incorporate trails exponential in number of recursive levels of trail applications

- **Problem:** (semi-)automatic trail creation → large number of uncertain trail definitions → large query rewrite time and large number of low quality results

- **Our Solution:** exploit uncertainty to improve rewrite time and precision
  - Prune rewrites by only using high-quality trails (top-K)
  - Prune rewrites by limiting trail recursivity (levels)
  - Prune rewrites by both (top-K, levels)

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# Effect of Trail Pruning: Performance

- Randomly generated trails that recurse with 1% chance
- Trail probabilities Zipf-distributed



Query-rewrite time controlled with pruning

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# Effect of Trail Pruning: Quality

- *Estimated* precision and number of relevant results returned using quality model



Pruning restricts attention to high quality trails only

Pruning reduces expected recall as less trails are used

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# Evaluation of MMCA with Pruning

- Randomly generated trails with mutual match chance of 1%
- Trail probabilities follow a Zipf distribution



Query-rewrite time controlled with pruning

Pruning restricts attention to high quality trails only

# Outline

2005      2006      2007      2008      2009

## 1. iMeMex Data Model (iDM)

Schema-agnostic representation of heterogeneous data (VLDB'06)

## 2. Semantic Trails

Pay-as-you-go Enrichment of a Dataspace through Relationships among Sets of Instances (VLDB'07)

## 3. Association Trails

Modeling of Fine-grained Relationships among Individual Instances in a Dataspace (in submission)

## 4. iMeMex PDSMS Architecture

System Architecture for Dataspace Management (VLDB'05 demo, CIDR'07 demo)

## 5. Lessons Learned

# Outline

2005 | 2006 | 2007 | 2008 | 2009

**1. iMeMex Data Model (iDM)**

Schema-agnostic representation of heterogeneous data (VLDB'06)

**2. Semantic Trails**

Pay-as-you-go Enrichment of a Dataspace through Relationships among Sets of Instances (VLDB'07)

**3. Association Trails**

Modeling of Fine-grained Relationships among Individual Instances in a Dataspace (in submission)

**4. iMeMex PDSMS Architecture**

System Architecture for Dataspace Management (VLDB'05 demo, CIDR'07 demo)

**5. Lessons Learned**

SAARLAND UNIVERSITY

COMPUTER SCIENCE

# Motivation: Social Networks Today



isFriend

Fred

isFriend

isFriend

isFriend

- **Query Services**
  - Keyword search
  - Browsing of friend lists and suggestions (e.g. same university)

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# Social Networks Tomorrow: An Overlay Graph

SAARLAND UNIVERSITY
COMPUTER SCIENCE

# Association Trails: Intensional Associations among Individual Instances

- **Our Approach**
  - Each association trail encodes a set of intensional edges in the association graph (e.g. <span style="color:red">sharesResearchInterest</span>)
  - Queries return not only primary results but also context in which they are in

- **Example:** when you query for "Fred", you would also get:
  - Comments people wrote about Fred
  - People who share research interests with Fred, or who went to the same university as Fred, or who graduated in the same year as Fred
  - All of the above, ranked by how much the item is related

# Definition of an Association Trail

- **Basic Form of an Association Trail**

Join Predicate that relates elements from the left with elements from the right

$$Q_L \quad \boxed{\theta(L, R)} \xrightarrow{\phantom{xx}} Q_R$$

- **Intuition:** When I return items from $Q_L$, you should also return related items from $Q_R$

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# Differences to Semantic Trails

| | Semantic Trails | Association Trails |
|---|---|---|
| **Type of Relationship** | Define equivalences among queries (sets) | Define relationships among instances |
| | | |
| | equivalent to //imap/marcos/iMeMex<br><br>• Query on //projects/PIM also returns //imap/marcos/iMeMex | hobbies are related<br><br>• Any query returning a person X also returns persons who share hobbies with X |

**Semantic Trails do not specify join semantics!**

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# Association Trail Examples

- **People who share research interests are related**

$$\text{sharesResearchInterest: //person} \overset{\theta_1}{\Longrightarrow} \text{//person,}$$

$$\theta_1(L,R) = (\exists i_1 \in L/\text{researchInterest:}$$

$$i_1 \in R/\text{researchInterest})$$

- **People who graduated in the same year are related**

$$\text{graduatedSameYear: //person} \underset{\theta_2}{\Longrightarrow} \text{//person,}$$

$$\theta_2(L,R) = (L.\text{gradYear} = R.\text{gradYear})$$

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# Answering Queries with Association Trails

- **Problem:** canonical plan to answer queries with association trails is expensive

- **Our Solutions**

    1. *N-Semi-Joins:* Reuse Common Subexpressions in trail queries

    2. *MatLR:* Index trail queries ($Q_L$ and $Q_R$)

    3. *MatFullJoin:* Materialize all of the intensional graph ($Q_L \bowtie Q_R$)

    4. *GCI:* Grouping-Compressed Index (join in linear space)

# Evaluation of Query Performance with Association Trails

- Synthetic social network data: 1.6 Million people
- Trail predicates are equi-joins on Zipf-distributed attributes



GCI offers more than an order of magnitude gain over Canonical

GCI offers more than an order of magnitude gain over MatFullJoin

# Outline



2005     2006     2007     2008     2009

**1. iMeMex Data Model (iDM)**

Schema-agnostic representation of heterogeneous data (VLDB'06)

**2. Semantic Trails**

Pay-as-you-go Enrichment of a Dataspace through Relationships among Sets of Instances (VLDB'07)

**3. Association Trails**

Modeling of Fine-grained Relationships among Individual Instances in a Dataspace (in submission)

**5. Lessons Learned**

**4. iMeMex PDSMS Architecture**

System Architecture for Dataspace Management (VLDB'05 demo, CIDR'07 demo)

SAARLAND UNIVERSITY
COMPUTER SCIENCE

# Outline

2005      2006      2007      2008      2009

**1. iMeMex Data Model (iDM)**

Schema-agnostic representation of heterogeneous data (VLDB'06)

**2. Semantic Trails**

Pay-as-you-go Enrichment of a Dataspace through Relationships among Sets of Instances (VLDB'07)

**3. Association Trails**

Modeling of Fine-grained Relationships among Individual Instances in a Dataspace (in submission)

**5. Lessons Learned**

**4. iMeMex PDSMS Architecture**

System Architecture for Dataspace Management (VLDB'05 demo, CIDR'07 demo)

# Design of the iMeMex PDSMS Architecture

- **Goals**
  - Hybrid architecture in-between search engines and information integration systems
  - *Not an information integration system:* no need to pre-declare schemas in order to query data
  - *Not a search engine:* allow for pay-as-you-go information integration
  - *Not a DBMS:* system does not take full control of the data
- **Our approach**
  - Use iQL to search and query all of the user's dataspace
  - Rewrite queries with semantic and association trails
  - Represent all the data in the sources with iDM

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# iMeMex Architecture: Logical Layers

- ## Logical Independence Layer (LIL)
  - Provides basic iQL planning
  - Handles all trail rewrites
- ## Physical Independence Layer (PHIL)
  - Abstracts from sources and formats, exposing a resource view graph

iQL queries

**iMeMex PDSMS**

Logical Independence Layer (LIL)

Physical Independence Layer (PHIL)

text, links

**Data Sources**

SAARLAND UNIVERSITY
COMPUTER SCIENCE

# iMeMex Architecture: Component Architecture



**iMeMex PDSMS**

**Logical Independence Layer (LIL): Query Processor**

- iQL Planner
- Trail Manager

**Physical Independence Layer (PHIL): Resource View Manager**

**Indexes & Replicas**
- Inverted Index
- Rowstore

**Content2iDM Converters**
- LaTeX
- XML
- MP3
- …

**Data Source Proxy**
- Filesystem
- IMAP
- Google
- …

Logical planning of iQL

Semantic and association trail rewrites

Indexing for data shipped to iMeMex

Convert formats to iDM

Data or query shipping

- Extensible: Everything is a plug-in (OSGi)
- Open-source (Apache 2.0): http://www.imemex.org

SAARLAND UNIVERSITY
COMPUTER SCIENCE

# iMeMex Prototype Statistics

- ~ 1500 classes

- ~ 127,000 LOC

- Java-based: supported on Linux, Mac and Windows

- OSGi-based: Everything is a Plug-in (~ 75 bundles)

- Open-source (Apache 2.0): http://www.imemex.org

# iMeMex Benefits: Views on your Desktop

- provides **unified concept** for handling unstructured, semi-structured and structured data on the user's desktop

- Allows users to define **arbitrary views**

- Queries can be specified using

  - keyword search
  - SQL
  - XQuery

```xml
<?xml version="1.0" encoding="utf-8" ?>
<imemex-query>
    <alias>
        <realname></realname>
        <name-in-query></name-in-query>
    </alias>
    <xquery><![CDATA[]]></xquery>
    <sql></sql>
    <keyword>test</keyword>
    <output-format></output-format>
</imemex-query>
```

# iMeMex Benefits: Views on your Desktop

- iMeMex QueryDispatcher plugin subscribes to *.query-files
- QueryDispatcher is responsible for executing that query

```xml
<?xml version="1.0" encoding="utf-8" ?>
<imemex-query>
    <alias>
        <realname></realname>
        <name-in-query></name-in-query>
    </alias>
    <xquery><![CDATA[]]></xquery>
    <sql></sql>
    <keyword>test</keyword>
    <output-format></output-format>
</imemex-query>
```

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# iMeMex Benefits: Views on your Desktop

- View results look like files/folders in the explorer
- The result of a view is only computed on demand! (when the user tries to read the content of a view)

virtual views vs. materialized views

# iMeMex Benefits: Views on your Desktop

```xml
1   <?xml version = '1.0' encoding = 'utf-8'?>
2   <imemex-query>
3       <alias>
4           <realname>file:///C:/tests/pim/ROOT/bla/student projects.xls</realname>
5           <name-in-query>file1</name-in-query>
6       </alias>
7       <alias>
8           <realname>file:///C:/tests/pim/ROOT/bla/student emails.doc</realname>
9           <name-in-query>file2</name-in-query>
10      </alias>
11      <xquery><![CDATA[
12          <result>{
13              for $a in doc("file1")//Student,
14              $b in doc("file2")//Name
15              where $a/text() = $b/text()
16              return   <person>
17                          <name> {$a/text()}</name>
18                          <projectType> {$a/../Type/text()} </projectType>
19                          <email> {$b/../Email/text()} </email>
20                      </person>
21          }</result>
22          ]]>
23      </xquery>
24      <sql/>
25      <search/>
26      <output-format>xls</output-format>
27  </imemex-query>
```

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# iMeMex Dataspace Navigator

# Related Systems Overview

- **Search Engines**
  - TopX [VLDB05], FleXPath [SIGMOD04], XSearch [VLDB03], XRank [SIGMOD03], Google [WWW98]

- **Information Integration Systems**
  - TSIMMIS (GAV) [ICDE95], Information Manifold (LAV) [VLDB96], GLAV [AAAI99], Piazza (P2P) [SIGMOD04], Multibase [VLDB83] , Garlic [VLDB97]

- **Dataspace Systems**
  - Dataspace vision [SIGMOD Record 05], Quarry [CIDR07, IIMAS08], PayGO [CIDR07]

- **PIM Systems**
  - MyLifeBits [SIGMOD05], Stuff I've Seen / Phlat [SIGIR03, CHI06], Haystack [CIDR05], SEMEX [CIDR05]

# Outline

2005      2006      2007      2008      2009

## 1. iMeMex Data Model (iDM)

Schema-agnostic representation of heterogeneous data (VLDB'06)

## 2. Semantic Trails

Pay-as-you-go Enrichment of a Dataspace through Relationships among Sets of Instances (VLDB'07)

## 3. Association Trails

Modeling of Fine-grained Relationships among Individual Instances in a Dataspace (in submission)

## 4. iMeMex PDSMS Architecture

System Architecture for Dataspace Management (VLDB'05 demo, CIDR'07 demo)

## 5. Lessons Learned

# Outline

2005     2006     2007     2008     2009

## 1. iMeMex Data Model (iDM)

Schema-agnostic representation of heterogeneous data (VLDB'06)

## 2. Semantic Trails

Pay-as-you-go Enrichment of a Dataspace through Relationships among Sets of Instances (VLDB'07)

## 3. Association Trails

Modeling of Fine-grained Relationships among Individual Instances in a Dataspace (in submission)

## 4. iMeMex PDSMS Architecture

System Architecture for Dataspace Management (VLDB'05 demo, CIDR'07 demo)

## 5. Lessons Learned

SAARLAND UNIVERSITY
COMPUTER SCIENCE

# Lessons Learned

2005      2006      2007      2008      2009

**1. iMeMex Data Model (iDM)**

Schema-agnostic representation of heterogeneous data (VLDB'06)

**Pros:**

abstracts away data formats

lean

flexible

yet powerful

clear separation: data model vs. data format

lazy features very useful e.g. simulating active xml

**Cons:**

no engine support (had to code everthing ourselves)

query processing on graphs requires effort

sometimes too powerful

lazy computation hard to control

sometimes too much OO-like

(our fault, not the model)

**What could be done in future:**

use RDF to implement iDM (scalable engines have only recently become available)

System Architecture for Dataspace Manageme

COMPUTER SCIENCE

# Lessons Learned

2005          2006          2007

## 1. iMeMex Data Model (iDM)

Schema-agnostic representation of heterogeneous data (VLDB'06)

## 2. Semantic Trails

Pay-as-you-go Enrichment of a Dataspace through Relationships among Sets of Instances (VLDB'07)

**What could be done in future:**

relevance feedback

automatic schema-matching

semi-automatic trail creation

application to Web

**Pros:**

easy

powerful

scalable

data model/format independent

intra-source relations

**Cons:**

applicability to fully structured data unclear

expressiveness of trails needs to be further investigated

=> impact on query rewrite?

SAARLAND UNIVERSITY
COMPUTER SCIENCE

# Lessons Learned

**Pros:**

easy

powerful

scalable

data model/format independent

several application domains: PIM, social networks, Web 3.0, ...

**Cons:**

indexing effort still high

**What could be done in future:**

eval on real social network

association advisor

update handling

trail sharing (as done today for delicious; wikipedia)

**3. Association Trails**

Modeling of Fine-grained Relationships among Individual Instances in a Dataspace (in submission)

200

**ils**

ent
ions
among Sets of Instances (VLDB'07)

**S Architecture**

ement (VLDB'05 demo, CIDR'07 demo)

5. Lessons

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# Lessons Learned

**Pros:**

OSGi

easy to prototype new functionality

hybrid mediation/ETL modes

full control

powerful

**Cons:**

OSGi

levels of abstraction hard to debug

a lot of functionality but sometimes unstable

testing difficult

**What could be done in future:**

do **not** use OSGi

offer less functionality, but get that right

better testing in the first place

„solve a smaller problem"
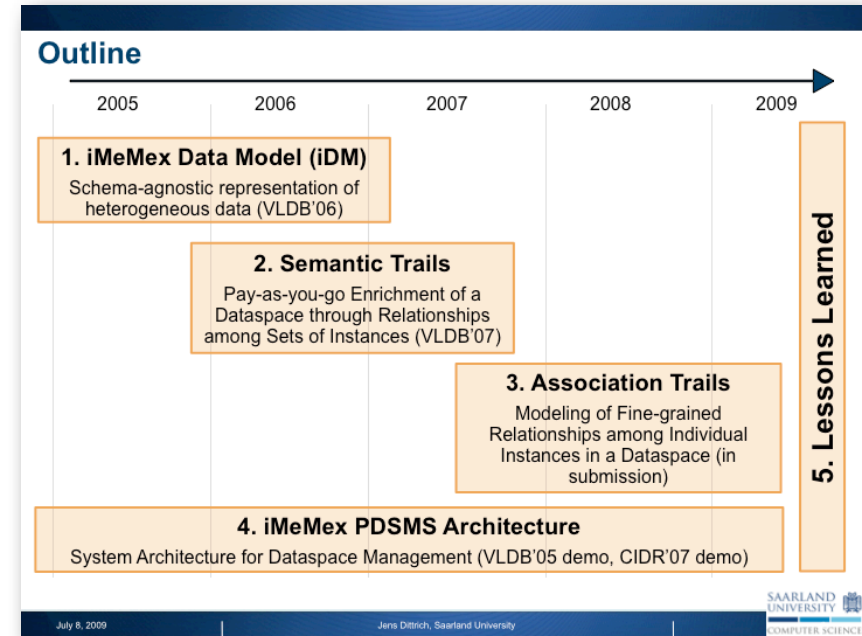
avoid code rot

**either** prototype **or** real system

dataspaces in the „cloud"

dataspace architectures

shared dataspaces

dataspaces and the Web

20

## 4. iMeMex PDSMS Architecture

System Architecture for Dataspace Management (VLDB'05 demo, CIDR'07 demo)

SAARLAND
UNIVERSITY
COMPUTER SCIENCE

# Conclusions on iMeMex

- ## Four contributions to the design of dataspace systems
  - ### iMeMex Data Model (iDM)
    - Logical Data Model for Personal Dataspaces
    - Lazily-computed Graph
  - ### iTrails
    - Pay-as-you-go Technique for Defining Relationships among Sets of Instances
  - ### Association Trails
    - Declarative Relationships among Individual Instances
  - ### iMeMex PDSMS Architecture
    - Architecture of a Personal Dataspace Management System



Outline

2005 | 2006 | 2007 | 2008 | 2009

**1. iMeMex Data Model (iDM)**
Schema-agnostic representation of heterogeneous data (VLDB'06)

**2. Semantic Trails**
Pay-as-you-go Enrichment of a Dataspace through Relationships among Sets of Instances (VLDB'07)

**3. Association Trails**
Modeling of Fine-grained Relationships among Individual Instances in a Dataspace (in submission)

**4. iMeMex PDSMS Architecture**
System Architecture for Dataspace Management (VLDB'05 demo, CIDR'07 demo)

**5. Lessons Learned**

July 8, 2009 — Jens Dittrich, Saarland University

SAARLAND UNIVERSITY
COMPUTER SCIENCE

# My personal Conclusions on Dataspaces

- great vision (still like it)
- at the same time too vague
- hard problems
- difficult to achieve
- did not become a hype (yet?)
- current buzz: "clouds", some overlap with dataspaces
- follow the crowd or do something risky
- maybe dataspace vision came 10 years too early
- some of the ideas behind "dataspaces" will re-appear under a different buzz word...

# Acknowledgments

- ETH Zurich
- SNF
- Donald Kossmann
- PhD students:
  - Marcos Salles
  - Lukas Blunschi
- MSC/BSc students:
  - Tobias Abt
  - Aarno Aukia
  - Sandro Blum

- Urs Blum
- Sibylle Dürr
- Markus Färber
- Steven Fluck
- Pascal Gamper
- Olivier Girard
- Stefan Hildenbrand
- Julia Imhof
- Roger Jäggi
- Shant Kirakos Karakashian
- André Schmidt

- Stefan Stalder
- Marco Steybe
- Philip Stutz
- Christian Tarnutzer
- Georg Troxler

Part of the slides in this talk from Marcos Salles: Thanks!

SAARLAND UNIVERSITY
COMPUTER SCIENCE