# XML Warehousing Meets Sociology

F.-X. Dudouet[1]    I. Manolescu[2]
B. Nguyen[3]    P. Senellart[2,4]

*IADIS ICWI*, October 20th 2005

**Introduction**
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

# Outline

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

# Sociological Process

1. Formulate hypotheses

2. Validate on data

   - Relevant sociological concepts (individuals, institutions. . . )
   - Data sources are: existing documents, interviews. . .

3. Conclude and issue new hypotheses

### Issue

How to collect and manage large volumes of heterogeneous information?

**Introduction**
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

# Sociological Process

1. Formulate hypotheses
2. Validate on data
   - Relevant sociological concepts (individuals, institutions. . . )
   - Data sources are: existing documents, interviews. . .
3. Conclude and issue new hypotheses

## Issue

How to collect and manage large volumes of heterogeneous information?

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

# Sociological Process

**1** Formulate hypotheses
**2** Validate on data
- Relevant sociological concepts (individuals, institutions...)
- Data sources are: existing documents, interviews...

**3** Conclude and issue new hypotheses

### Issue

How to collect and manage large volumes of heterogeneous information?

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

# Sociological Process

1. Formulate hypotheses
2. Validate on data
   - Relevant sociological concepts (individuals, institutions. . . )
   - Data sources are: existing documents, interviews. . .
3. Conclude and issue new hypotheses

## Issue

How to collect and manage large volumes of heterogeneous information?

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

# Sociological Process

1. Formulate hypotheses
2. Validate on data
   - Relevant sociological concepts (individuals, institutions. . . )
   - Data sources are: existing documents, interviews. . .
3. Conclude and issue new hypotheses

## Issue

How to collect and manage large volumes of heterogeneous information?

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

# Sociological Process

1. Formulate hypotheses
2. Validate on data
   - Relevant sociological concepts (individuals, institutions. . . )
   - Data sources are: existing documents, interviews. . .
3. Conclude and issue new hypotheses

## Issue

How to collect and manage large volumes of heterogeneous information?

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

## Case of the World Wide Web

- Inestimable source of data
- Much human activity involve Web technology

But:

- Heterogeneity of sources
- Not suited to classical database systems
- Need of conceptual models

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

# Case of the World Wide Web

- Inestimable source of data
- Much human activity involve Web technology

But:

- Heterogeneity of sources
- Not suited to classical database systems
- Need of conceptual models

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

# Case of the World Wide Web

- Inestimable source of data
- Much human activity involve Web technology

But:

- Heterogeneity of sources
- Not suited to classical database systems
- Need of conceptual models

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

# Case of the World Wide Web

- Inestimable source of data
- Much human activity involve Web technology

But:

- Heterogeneity of sources
- Not suited to classical database systems
- Need of conceptual models

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

# Case of the World Wide Web

- Inestimable source of data
- Much human activity involve Web technology

But:

- Heterogeneity of sources
- Not suited to classical database systems
- Need of conceptual models

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

# Standardization

## Standard negociations

$\implies$ Important economic and political impact

## Issue

Who? Why? How?

## Example

XQuery standardization scene

- Arena quite accessible via mailing lists

- Author's acquaintance with the topic

- Process almost finished

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
**Standardization**

# Standardization

Standard negociations

$\implies$ Important economic and political impact

## Issue

Who? Why? How?

## Example

XQuery standardization scene

- Arena quite accessible via mailing lists
- Author's acquaintance with the topic
- Process almost finished

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

# Standardization

Standard negociations

$\implies$ Important economic and political impact

## Issue

Who? Why? How?

## Example

XQuery standardization scene

- Arena quite accessible via mailing lists

- Author's acquaintance with the topic

- Process almost finished

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

# Standardization

Standard negociations
$\implies$ Important economic and political impact

## Issue

Who? Why? How?

## Example

XQuery standardization scene

- Arena quite accessible via mailing lists
- Author's acquaintance with the topic
- Process almost finished

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
**Standardization**

# Standardization

Standard negociations
$\implies$ Important economic and political impact

## Issue

Who? Why? How?

## Example

XQuery standardization scene

- Arena quite accessible via mailing lists
- Author's acquaintance with the topic
- Process almost finished

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
**Standardization**

# Standardization

Standard negociations

$\implies$ Important economic and political impact

## Issue

Who? Why? How?

## Example

XQuery standardization scene

- Arena quite accessible via mailing lists
- Author's acquaintance with the topic
- Process almost finished

Introduction
Methodology
Experimentation
Conclusion

Sociological Process
Standardization

# Standardization

Standard negociations
$\implies$ Important economic and political impact

## Issue

Who? Why? How?

## Example

XQuery standardization scene

- Arena quite accessible via mailing lists
- Author's acquaintance with the topic
- Process almost finished

Introduction
Methodology
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

# Outline

Introduction
**Methodology**
Experimentation
Conclusion

**Conceptual process**
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

# Modelling and analysis process

- Modelling the relevant sociological entities (actors, institutions, functions, messages, time)
- Designing a warehouse of Web resources relevant to the sociological analysis
- Exploiting the warehouse (feeding the warehouse, issuing queries)

Queries enable verification of the hypotheses

Introduction
**Methodology**
Experimentation
Conclusion

**Conceptual process**
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

# Modelling and analysis process

- Modelling the relevant sociological entities (actors, institutions, functions, messages, time)
- Designing a warehouse of Web resources relevant to the sociological analysis
- Exploiting the warehouse (feeding the warehouse, issuing queries)

Queries enable verification of the hypotheses

Introduction
**Methodology**
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

## Modelling and analysis process

- Modelling the relevant sociological entities (actors, institutions, functions, messages, time)
- Designing a warehouse of Web resources relevant to the sociological analysis
- Exploiting the warehouse (feeding the warehouse, issuing queries)

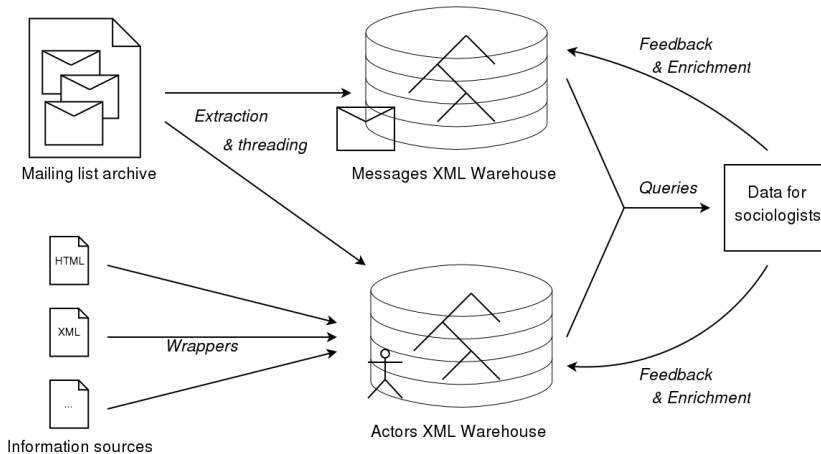Queries enable verification of the hypotheses

Introduction
**Methodology**
Experimentation
Conclusion

**Conceptual process**
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

## Modelling and analysis process

- Modelling the relevant sociological entities (actors, institutions, functions, messages, time)
- Designing a warehouse of Web resources relevant to the sociological analysis
- Exploiting the warehouse (feeding the warehouse, issuing queries)

Queries enable verification of the hypotheses

Introduction
Methodology
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

# Warehouse construction process

Introduction
Methodology
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

# XML Warehousing

## Pros

- **Semi-structured** information
- Tree structure of a mailing list
- Simple to understand

Queries on XML warehouses: XQuery itself!

Introduction
Methodology
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

# XML Warehousing

## Pros

- **Semi-structured** information
- **Tree structure** of a mailing list
- Simple to understand

Queries on XML warehouses: XQuery itself!

Introduction
Methodology
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

# XML Warehousing

## Pros

- **Semi-structured** information
- **Tree structure** of a mailing list
- **Simple** to understand

Queries on XML warehouses: XQuery itself!

Introduction
**Methodology**
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

# XML Warehousing

## Pros

- **Semi-structured** information
- **Tree structure** of a mailing list
- **Simple** to understand

Queries on XML warehouses: **XQuery** itself!

Introduction
Methodology
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

# Data filtering and enrichment

- Identify real-world objects represented in the warehouse
  - First name, last name, institution from e-mails
  - Identifying institutions participating in the process

- Classify these objects according to application-driven criteria
  - Issue classification queries to populate interesting classes (iterative process)

Introduction
**Methodology**
Experimentation
Conclusion

Conceptual process
XML Warehousing
**Data filtering and enrichment**
Complementary sociological tools

# Data filtering and enrichment

- Identify real-world objects represented in the warehouse
  - First name, last name, institution from e-mails
  - Identifying institutions participating in the process

- Classify these objects according to application-driven criteria
  - Issue classification queries to populate interesting classes (iterative process)

Introduction
Methodology
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

# Data filtering and enrichment

- Identify real-world objects represented in the warehouse
    - First name, last name, institution from e-mails
    - Identifying institutions participating in the process

- Classify these objects according to application-driven criteria
    Issue classification queries to populate interesting classes (iterative process)

Introduction
Methodology
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

# Data filtering and enrichment

- Identify real-world objects represented in the warehouse
  - First name, last name, institution from e-mails
  - Identifying institutions participating in the process

- Classify these objects according to application-driven criteria
    Issue classification queries to populate interesting classes (iterative process)

Introduction
**Methodology**
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
**Complementary sociological tools**

# Complementary sociological tools

## Issue

### Information on the Web has holes

- **Missing** information
- Important dimensions (e.g. time) implicitly or not at all represented
- Need to cross various sources to establish information

## Tools

- Interviews, inside information
- Human-readable data sources
- Statistics tools (social properties and group extraction)
- Human annotation and validation

Introduction
**Methodology**
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
**Complementary sociological tools**

# Complementary sociological tools

## Issue

Information on the Web has holes

- Missing information
- Important dimensions (e.g. time) implicitly or not at all represented
- Need to cross various sources to establish information

## Tools

- Interviews, inside information
- Human-readable data sources
- Statistics tools (social properties and group extraction)
- Human annotation and validation

Introduction
Methodology
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

# Complementary sociological tools

## Issue

Information on the Web has holes

- Missing information
- Important dimensions (e.g. time) implicitly or not at all represented
- Need to cross various sources to establish information

## Tools

- Interviews, inside information
- Human-readable data sources
- Statistics tools (social properties and group extraction)
- Human annotation and validation

Introduction
Methodology
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

# Complementary sociological tools

## Issue

Information on the Web has holes

- Missing information
- Important dimensions (e.g. time) implicitly or not at all represented
- Need to cross various sources to establish information

## Tools

- Interviews, inside information
- Human-readable data sources
- Statistics tools (social properties and group extraction)
- Human annotation and validation

Introduction
**Methodology**
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
**Complementary sociological tools**

# Complementary sociological tools

## Issue

Information on the Web has holes

- Missing information
- Important dimensions (e.g. time) implicitly or not at all represented
- Need to cross various sources to establish information

## Tools

- Interviews, inside information
- Human-readable data sources
- Statistics tools (social properties and group extraction)
- Human annotation and validation

Introduction
Methodology
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
Complementary sociological tools

# Complementary sociological tools

## Issue

Information on the Web has holes

- Missing information
- Important dimensions (e.g. time) implicitly or not at all represented
- Need to cross various sources to establish information

## Tools

- Interviews, inside information
- Human-readable data sources
- Statistics tools (social properties and group extraction)
- Human annotation and validation

Introduction
**Methodology**
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
**Complementary sociological tools**

# Complementary sociological tools

## Issue

Information on the Web has holes

- Missing information
- Important dimensions (e.g. time) implicitly or not at all represented
- Need to cross various sources to establish information

## Tools

- Interviews, inside information
- Human-readable data sources
- Statistics tools (social properties and group extraction)
- Human annotation and validation

Introduction
**Methodology**
Experimentation
Conclusion

Conceptual process
XML Warehousing
Data filtering and enrichment
**Complementary sociological tools**

# Complementary sociological tools

## Issue

Information on the Web has holes

- Missing information
- Important dimensions (e.g. time) implicitly or not at all represented
- Need to cross various sources to establish information

## Tools

- Interviews, inside information
- Human-readable data sources
- Statistics tools (social properties and group extraction)
- Human annotation and validation

Introduction
Methodology
**Experimentation**
Conclusion

Warehouses
Results
Sociological interpretation

# Outline

Introduction
Methodology
**Experimentation**
Conclusion

**Warehouses**
Results
Sociological interpretation

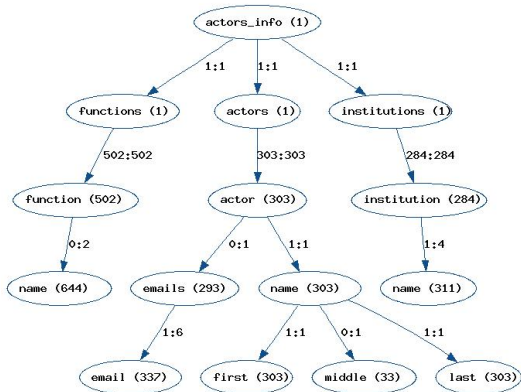## Message warehouse

`public-qt-comments@w3.org` mailing list.

### Data

- 5626 messages
- 2718 threads
- Maximum thread depth: 12

Introduction
Methodology
**Experimentation**
Conclusion

**Warehouses**
Results
Sociological interpretation

# Message warehouse

`public-qt-comments@w3.org` mailing list.

## Data

- 5626 messages
- 2718 threads
- Maximum thread depth: 12

Introduction
Methodology
**Experimentation**
Conclusion

**Warehouses**
Results
Sociological interpretation

# Message warehouse

`public-qt-comments@w3.org` mailing list.

## Data

- 5626 messages
- 2718 threads
- Maximum thread depth: 12

Introduction
Methodology
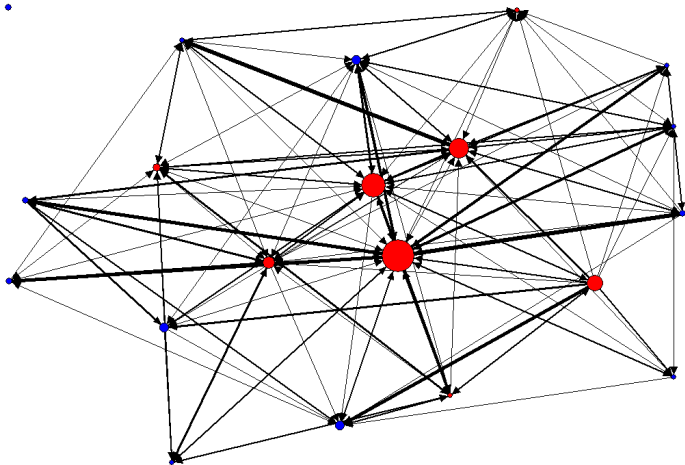**Experimentation**
Conclusion

**Warehouses**
Results
Sociological interpretation

# Message warehouse

`public-qt-comments@w3.org` mailing list.

## Data

- 5626 messages
- 2718 threads
- Maximum thread depth: 12

Introduction
Methodology
Experimentation
Conclusion

Warehouses
Results
Sociological interpretation

# Actors warehouse

Introduction
Methodology
**Experimentation**
Conclusion

Warehouses
Results
Sociological interpretation

# Simple results

Actor repartition and volume of interaction by affiliation profile

| Profile | # actors | # messages |
|---|---|---|
| Companies | 135 | 2689 |
| Universities | 39 | 112 |
| Organizations | 33 | 197 |
| Companies & Universities | 3 | 532 |
| Companies & Organizations | 22 | 1052 |
| Universities & Organizations | 6 | 36 |
| Non specified | 65 | 681 |
| **Total** | **303** | **5299** |

Introduction
Methodology
**Experimentation**
Conclusion

Warehouses
Results
Sociological interpretation

# Answer network

Introduction
Methodology
**Experimentation**
Conclusion

Warehouses
Results
Sociological interpretation

# Sociological interpretation

- Companies dominate XQuery standardization
- Key actors tend to have multiple affiliation
- Different profiles of participation in the list, even for key actors.

Introduction
Methodology
**Experimentation**
Conclusion

Warehouses
Results
Sociological interpretation

# Sociological interpretation

- Companies dominate XQuery standardization
- Key actors tend to have multiple affiliation
- Different profiles of participation in the list, even for key actors.

Introduction
Methodology
**Experimentation**
Conclusion

Warehouses
Results
Sociological interpretation

# Sociological interpretation

- Companies dominate XQuery standardization
- Key actors tend to have multiple affiliation
- Different profiles of participation in the list, even for key actors.

Introduction
Methodology
Experimentation
**Conclusion**

Summary
Perspectives

# Outline

Introduction
Methodology
Experimentation
**Conclusion**

**Summary**
Perspectives

# Summary

- **Interdisciplinary** approach
- Use of semi-structured technology for sociological study
- Built an XML warehouse based on XQuery public W3C information
- Preliminary analysis of the warehouse data
- Companies seem to be first in standardization process

Introduction
Methodology
Experimentation
Conclusion

Summary
Perspectives

# Summary

- **Interdisciplinary** approach
- Use of **semi-structured** technology for **sociological** study
- Built an XML warehouse based on XQuery public W3C information
- Preliminary analysis of the warehouse data
- Companies seem to be first in standardization process

Introduction
Methodology
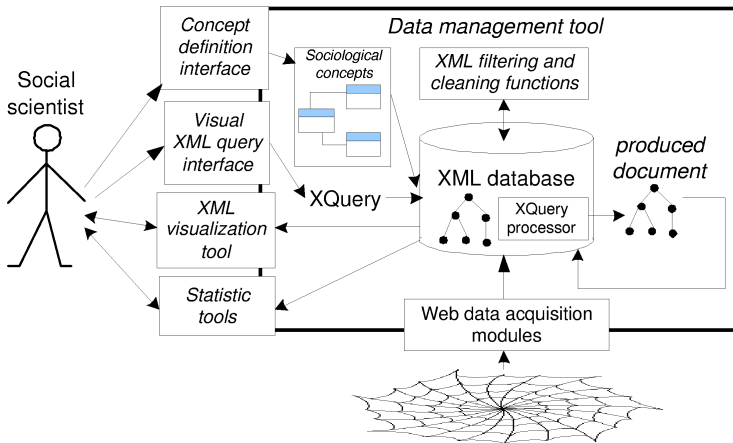Experimentation
Conclusion

Summary
Perspectives

# Summary

- **Interdisciplinary** approach
- Use of **semi-structured** technology for **sociological** study
- Built an **XML warehouse** based on XQuery public W3C information
- **Preliminary analysis** of the warehouse data
- Companies seem to be **first in standardization process**

Introduction
Methodology
Experimentation
Conclusion

Summary
Perspectives

# Summary

- Interdisciplinary approach
- Use of semi-structured technology for sociological study
- Built an XML warehouse based on XQuery public W3C information
- Preliminary analysis of the warehouse data
- Companies seem to be first in standardization process

Introduction
Methodology
Experimentation
Conclusion

Summary
Perspectives

# Summary

- **Interdisciplinary** approach
- Use of **semi-structured** technology for **sociological** study
- Built an **XML warehouse** based on XQuery public W3C information
- **Preliminary analysis** of the warehouse data
- Companies seem to be **first in standardization process**

Introduction
Methodology
Experimentation
**Conclusion**

Summary
**Perspectives**

# Generic Framework for the Social Scientist

Introduction
Methodology
Experimentation
**Conclusion**

Summary
**Perspectives**

# Future Work

- Textual analysis of message contents (e.g. agree/disagree)
- Proper management of temporal dimension
- Enriched actor warehouse with more sources (WWW in particular)
- Similar work on larger/other/private mailing lists
- More complex queries

Introduction
Methodology
Experimentation
Conclusion

Summary
Perspectives

# Future Work

- Textual analysis of message contents (e.g. agree/disagree)
- Proper management of temporal dimension
- Enriched actor warehouse with more sources (WWW in particular)
- Similar work on larger/other/private mailing lists
- More complex queries

Introduction
Methodology
Experimentation
Conclusion

Summary
Perspectives

# Future Work

- Textual analysis of message contents (e.g. agree/disagree)
- Proper management of temporal dimension
- Enriched actor warehouse with more sources (WWW in particular)
- Similar work on larger/other/private mailing lists
- More complex queries

Introduction
Methodology
Experimentation
**Conclusion**

Summary
Perspectives

# Future Work

- Textual analysis of message contents (e.g. agree/disagree)
- Proper management of temporal dimension
- Enriched actor warehouse with more sources (WWW in particular)
- Similar work on larger/other/private mailing lists
- More complex queries

Introduction
Methodology
Experimentation
Conclusion

Summary
Perspectives

# Future Work

- Textual analysis of message contents (e.g. agree/disagree)
- Proper management of temporal dimension
- Enriched actor warehouse with more sources (WWW in particular)
- Similar work on larger/other/private mailing lists
- More complex queries