# An unsupervised classification process for large datasets based on web reasoning

Rafael PEIXOTO, Thomas HASSAN, Christophe CRUZ, Aurelie BERTAUX, Nuno SILVA
thomas.hassan@u-bourgogne.fr

Laboratoire LE2I – UMR CNRS 6306 – Université de Bourgogne

**Context**
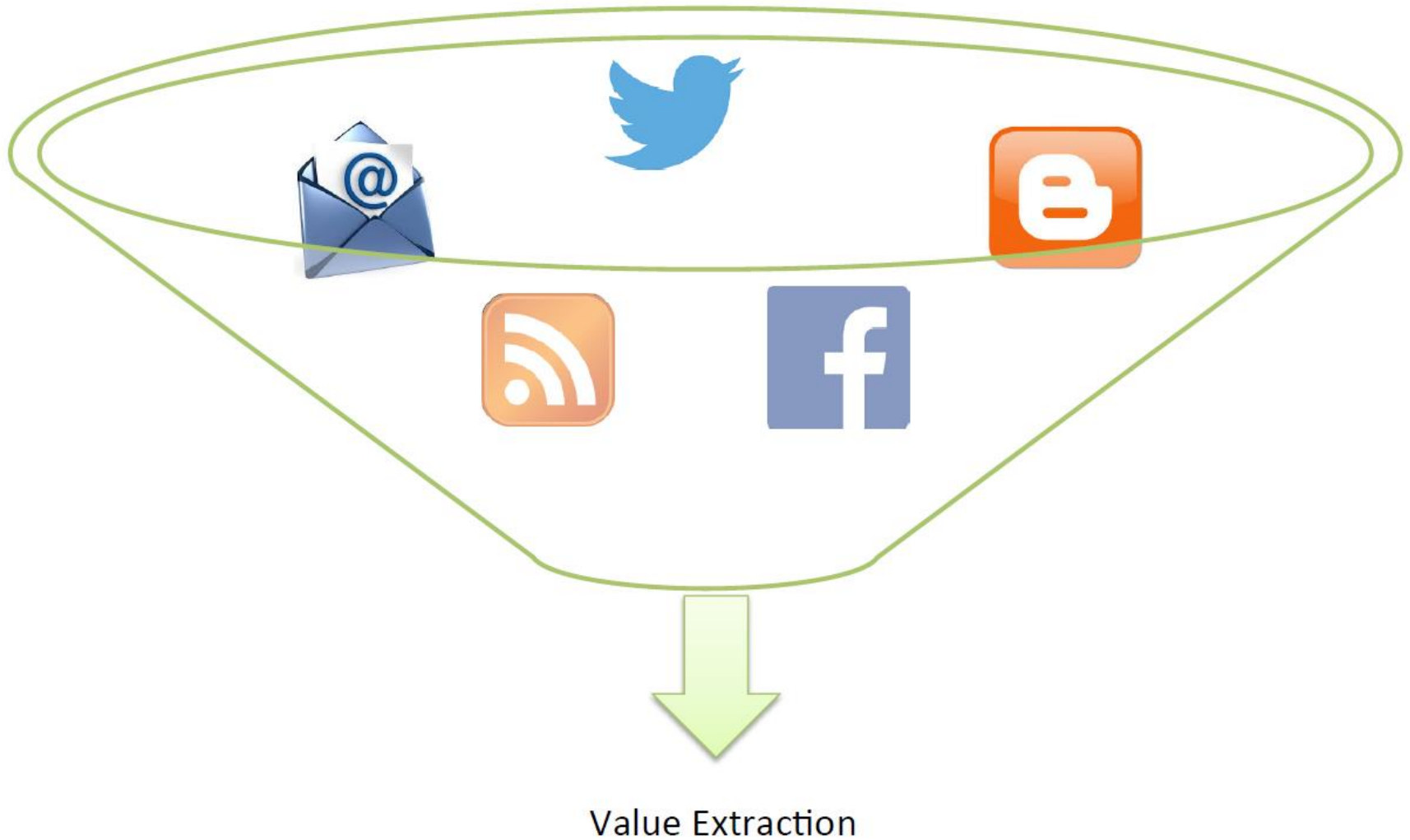
- Global problem
- The Semantic HMC

Specific Problem
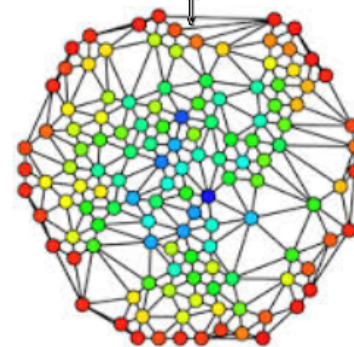
- Proposed Solution

Implementation

- Setup
- Results
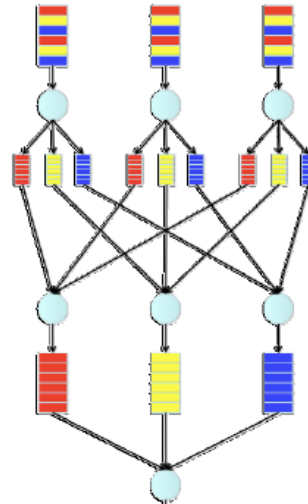
Conclusion and future work

**Value extraction** from **Big Data** sources



Value Extraction

○ **Why ontologies**

- Ontologies are the most accepted way to represent semantics in the Semantic Web and a good solution for intelligent computer systems that operate close to the human concept level, bridging the gap between human conceptions and computational requirements.

○ **Semantic HMC**

- Ontology-described predictive model

- Learned using Big Data technologies
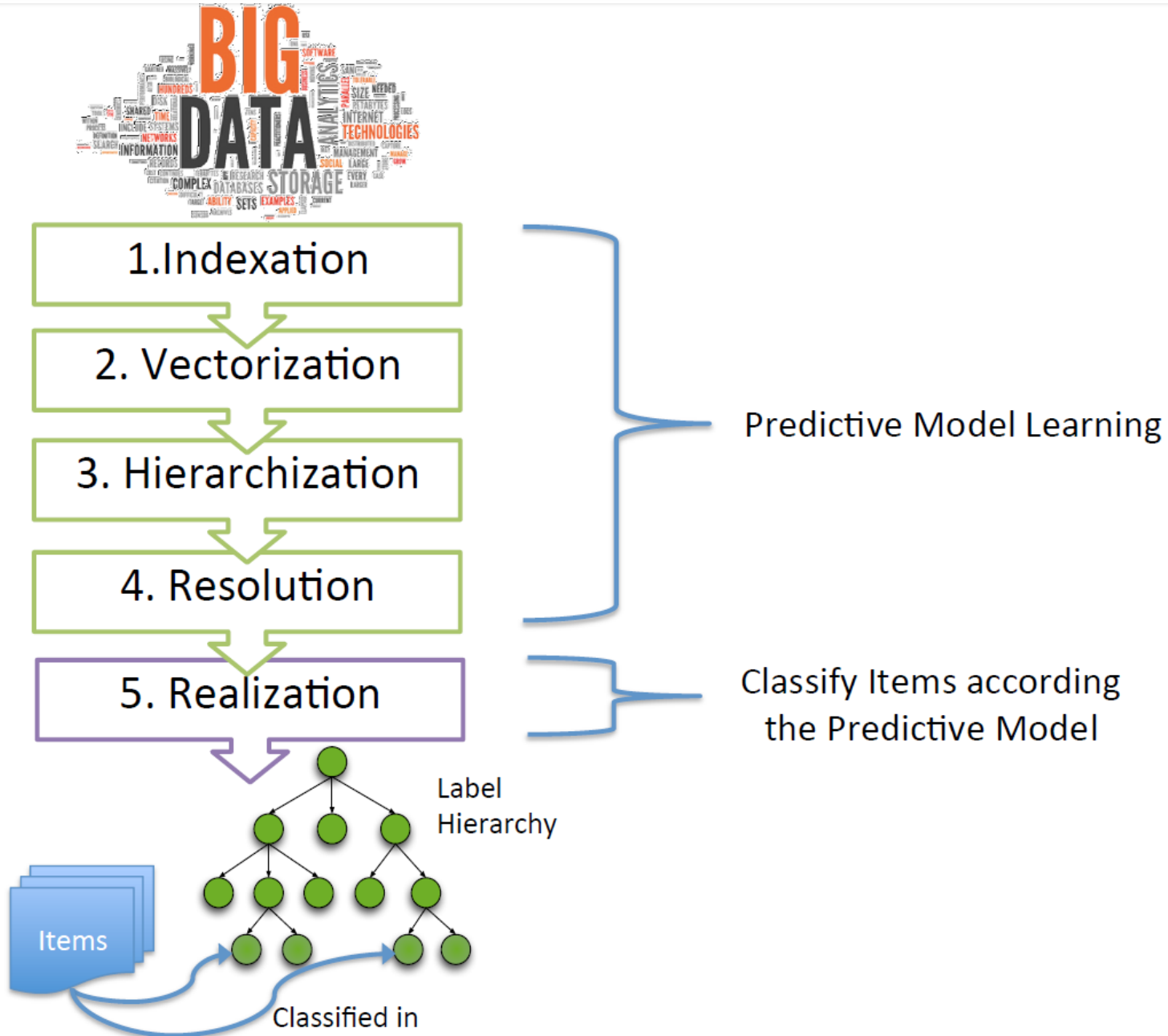
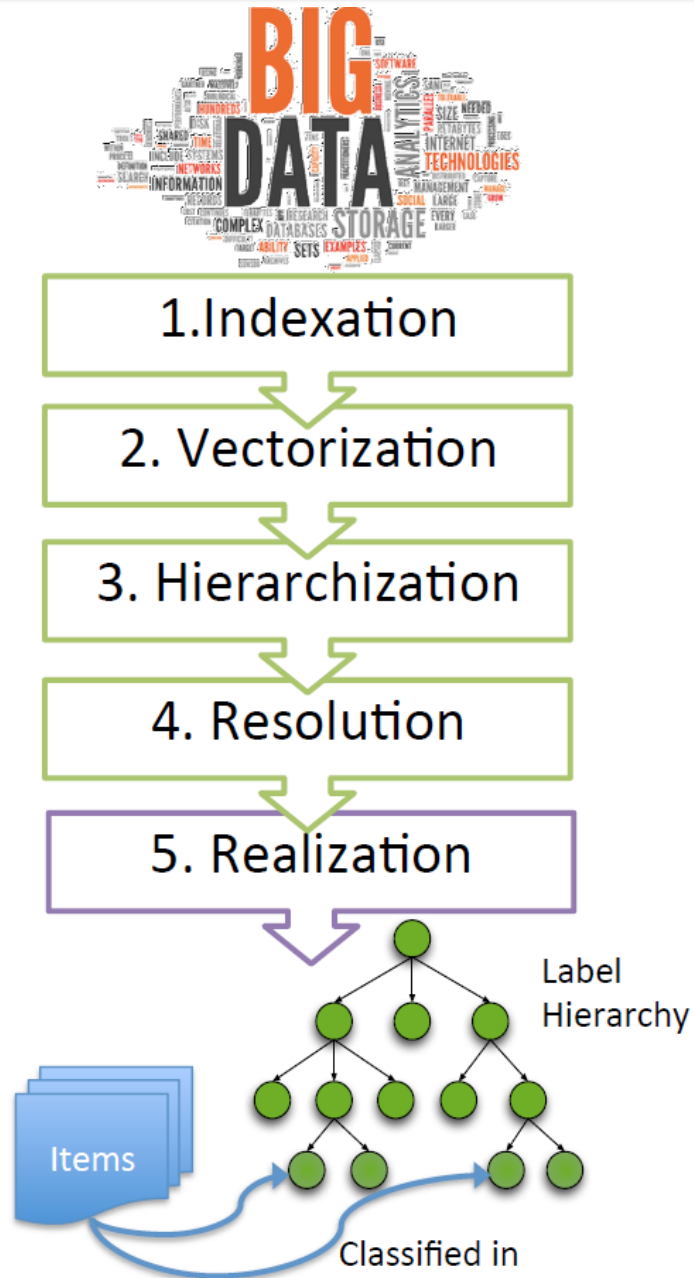- Rule-based Web Reasoning to perform classification
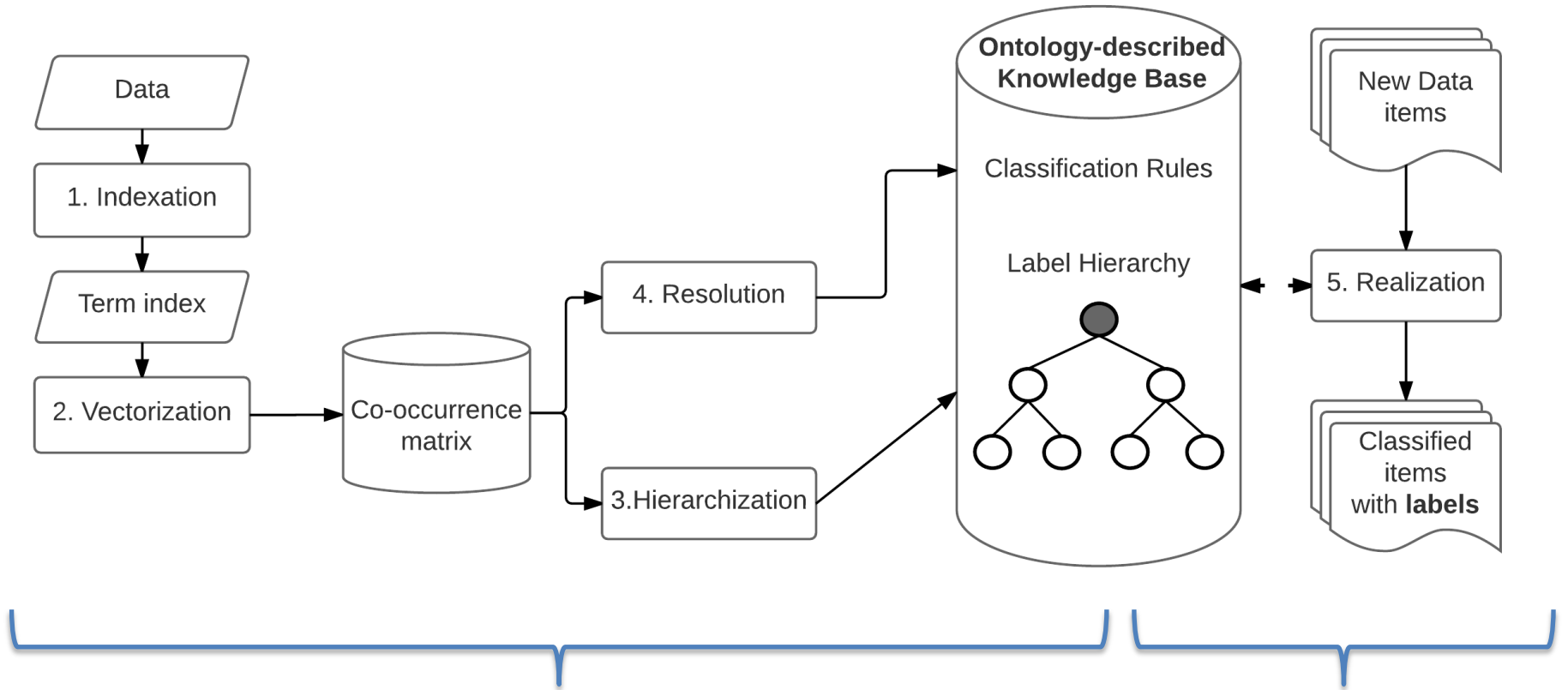
Big Data Technologies

S-HMC Processing

1.Indexation

2. Vectorization

3. Hierarchization

4. Resolution

5. Realization

Items

Label Hierarchy

Classified in

- o Indexation
  - o Extract terms
  - o Index the items
- o Vectorization
  - o Calculate term frequency vectors
  - o Co-occurrence matrix
- o Hierarchization
  - o Label selection
  - o Hierarchical relations
- o Resolution
  - o Classification rules creation
- o Realization
  - o Ontology population
  - o Rule-based Web Reasoning to classify items

7

Unsupervised **ontology learning**

**Rule-based**
Classification
(Web Reasoner)

Context

- Global problem
- The Semantic HMC
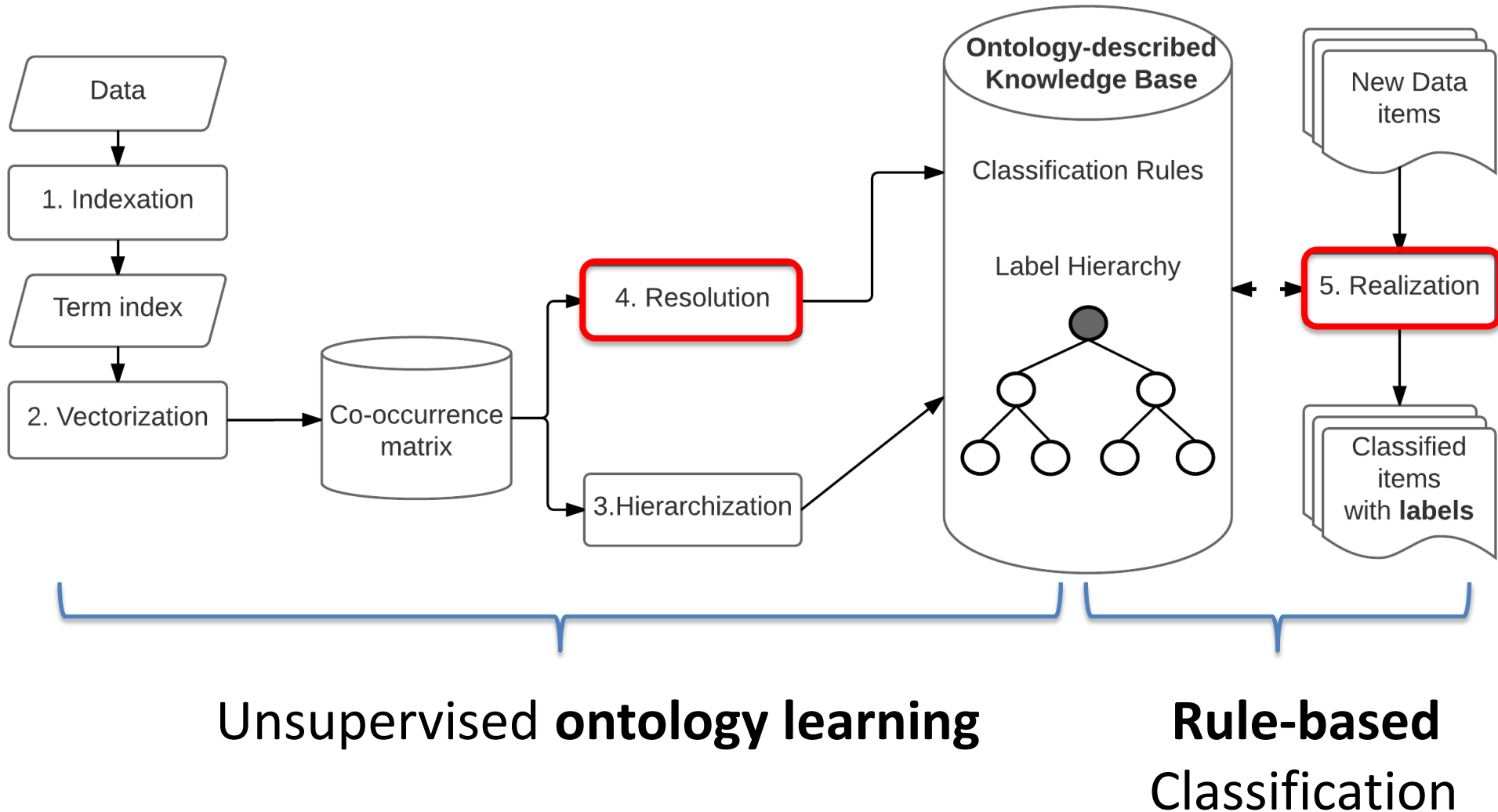
**Specific Problem**

- Proposed Solution

Implementation

- Setup
- Results

Conclusion and future work

**Rule-based** reasonning to perform **Classification**



Unsupervised **ontology learning**          **Rule-based** Classification

- *Resolution:* Learn **classifications rules** from **large volumes** of unstructured text

  ➡ Distributed method that exploits the coocurrence matrix

- *Realization:* classify **large volumes** of new **data items**
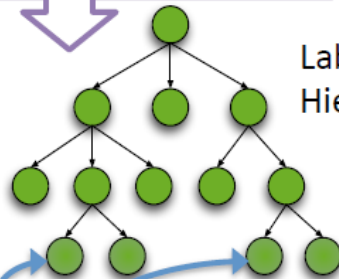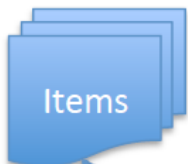
  ➡ Classification using a Web Reasonner

## Learning **Alpha** and **Beta** sets

| $P_C$(i\|j) | term$_1$ | term$_2$ | term$_3$ | term$_4$ | term$_5$ | term$_6$ | term$_7$ |
|---|---|---|---|---|---|---|---|
| label$_1$ | 0 | 0 | 5 | 0 | 5 | 25 | 25 |
| label$_2$ | 0 | 75 | 0 | 0 | 0 | 75 | 5 |
| label$_3$ | 0 | 0 | 75 | 0 | 25 | 0 | 0 |
| label$_4$ | 5 | 25 | 25 | 0 | 5 | 93 | 25 |
| label$_5$ | 95 | 0 | 0 | 0 | 60 | 0 | 5 |
| label$_6$ | 0 | 60 | 0 | 95 | 0 | 0 | 90 |
| label$_7$ | 5 | 98 | 5 | 60 | 25 | 0 | 79 |

**Coocurrence:**
$$P_C(term_i|term_j) = \frac{cfm\,(term_i, term_j)}{cfm(term_j, term_j)}$$

**Alpha set:**
$$\omega_\alpha^{t_i} = \{t_j | \forall t_j \in Term\colon P_C(t_i|t_j) > \alpha\}$$

**Beta set:**
$$\omega_\beta^{t_i} = \{t_j | \forall t_j \in Term\colon \beta \le P_C(t_i|t_j) \le \alpha\}$$

## Learning **Alpha** and **Beta** sets



**Alpha set:** $\quad \omega_\alpha^{t_i} = \{t_j | \forall t_j \in Term : P_C(t_i|t_j) > \alpha\}$

**Beta set:** $\quad \omega_\beta^{t_i} = \{t_j | \forall t_j \in Term : \beta \leq P_C(t_i|t_j) \leq \alpha\}$

**Example:**

| % | term$_1$ | term$_2$ | term$_3$ | term$_4$ | term$_5$ | term$_6$ | term$_7$ |
|---|---|---|---|---|---|---|---|
| label$_1$ | 0 | 0 | 5 | 0 | 5 | 25 | 25 |
| label$_2$ | 0 | 75 | 0 | 0 | 0 | 75 | 5 |
| label$_3$ | 0 | 0 | 75 | 0 | 25 | 0 | 0 |
| label$_4$ | 5 | 25 | 25 | 0 | 5 | 93 | 25 |
| label$_5$ | 95 | 0 | 0 | 0 | 60 | 0 | 5 |
| label$_6$ | 0 | 60 | 0 | 95 | 0 | 0 | 90 |
| label$_7$ | 5 | 98 | 5 | 60 | 25 | 0 | 79 |

$$\omega_\alpha^{t_i} = \{t_j | \forall t_j \in Term: P_C(t_i|t_j) > \alpha\}, \alpha = \boxed{91}$$

$$\omega_\beta^{t_i} = \{t_j | \forall t_j \in Term: \beta \leq P_C(t_i|t_j) \leq \alpha\}, \beta = \boxed{70}$$

## Classification at **query-time** using **backward-chaining**

# Core Ontology

| DL concepts | Description |
|---|---|
| $Item \sqsubseteq \exists hasTerm.Term$ | Items to classify (e.g. doc) has terms |
| $Term \sqsubseteq \top$ | Terms (e.g. word) extracted from items |
| $Label \sqsubseteq Term$ | Labels are terms used to classify items |
| $Label \sqsubseteq \forall broader.Label$ | Broader relation between labels |
| $Label \sqsubseteq \forall narrower.Label$ | Narrower relation between labels |
| $broader \equiv narrower^-$ | Broader and narrower are inverse |
| $Item \sqcap Term = \emptyset$ | Items and Terms are disjoint |
| $Item \sqsubseteq \forall isClassified.Label$ | Relation that links items with labels |

Context
- Global problem
- The Semantic HMC

Specific Problem
- Proposed Solution

**Implementation**
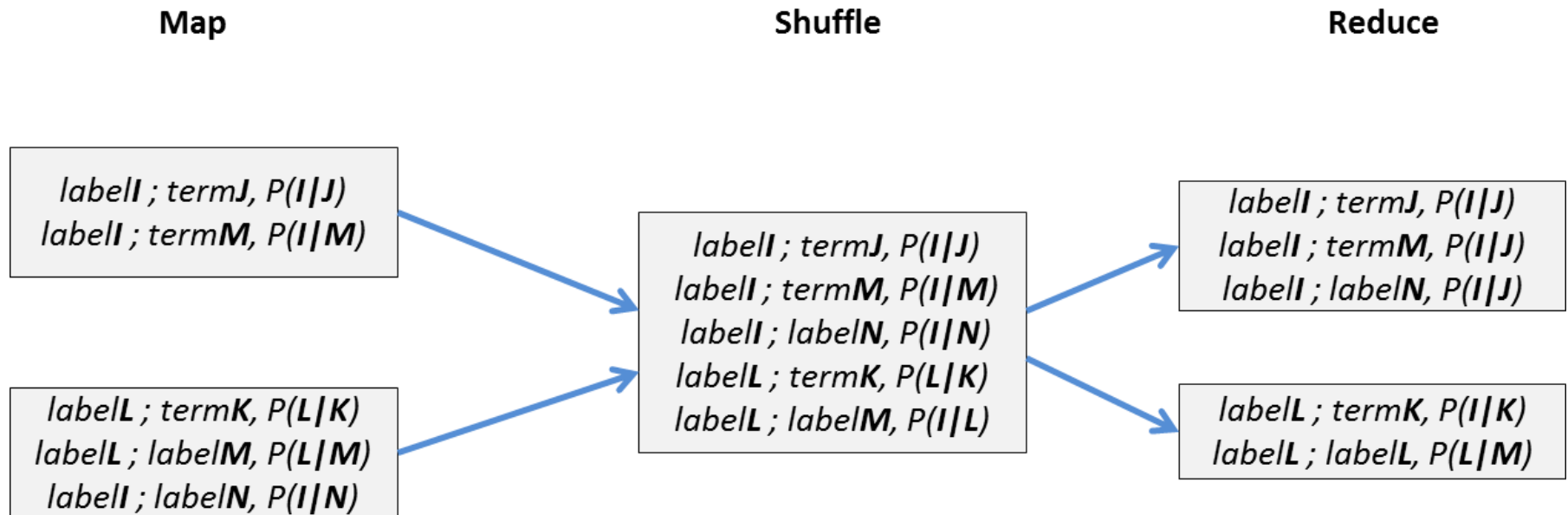- Setup
- Results

Conclusion and future work

**Distributed** process using mapreduce:



OWL API used to generate SWRL rules from the output

$$Item(?it), Term(term_i), Label(term_i), hasTerm\left(?it, term_j\right) \rightarrow isClassified(?it, term_i)$$

## Generated rules Exemple

| Alpha rules |
|---|
| $Item(?\,it), Term(t_1), Label(term_i), hasTerm(?\,it, t_1) \rightarrow$ $isClassified(?\,it, term_i)$ |
| $Item(?\,it), Term(t_2), Label(term_i), hasTerm(?\,it, t_2) \rightarrow$ $isClassified(?\,it, term_i)$ |

| Beta rules |
|---|
| $Item(?\,it), Term(t_1), Term(t_2), Label(term_i),$ $hasTerm(?\,it, t_1), hasTerm(?\,it, t_2) \rightarrow isClassified(?\,it, term_i)$ |
| $Item(?\,it), Term(t_1), Term(t_3), Label(term_i),$ $hasTerm(?\,it, t_1), hasTerm(?\,it, t_3) \rightarrow isClassified(?\,it, term_i)$ |
| $Item(?\,it), Term(t_2), Term(t_3), Label(term_i),$ $hasTerm(?\,it, t_2), hasTerm(?\,it, t_3) \rightarrow isClassified(?\,it, term_i)$ |

**Stardog** used as a scalable triple-store (compatible with **backward-chaining** inference as well as **SWRL** rules inference)

Rule selection process developped in Java interacting with Stardog to optimize query performance

## Dataset



WIKIPEDIA
*The Free Encyclopedia*

| Sub-Dataset | Number of articles |
| --- | --- |
| Wikipedia 1 | 174900 |
| Wikipedia 2 | 407000 |
| Wikipedia 3 | 994000 |

## Cluster



Google Cloud Platform

| Resource type | Description |
| --- | --- |
| Number of nodes | 4 |
| CPU (per node) | Intel Xeon E5 v2 |
| RAM (per node) | 7.5GB |
| Disk (per node) | 500GB |

# Implementation: parameter setup

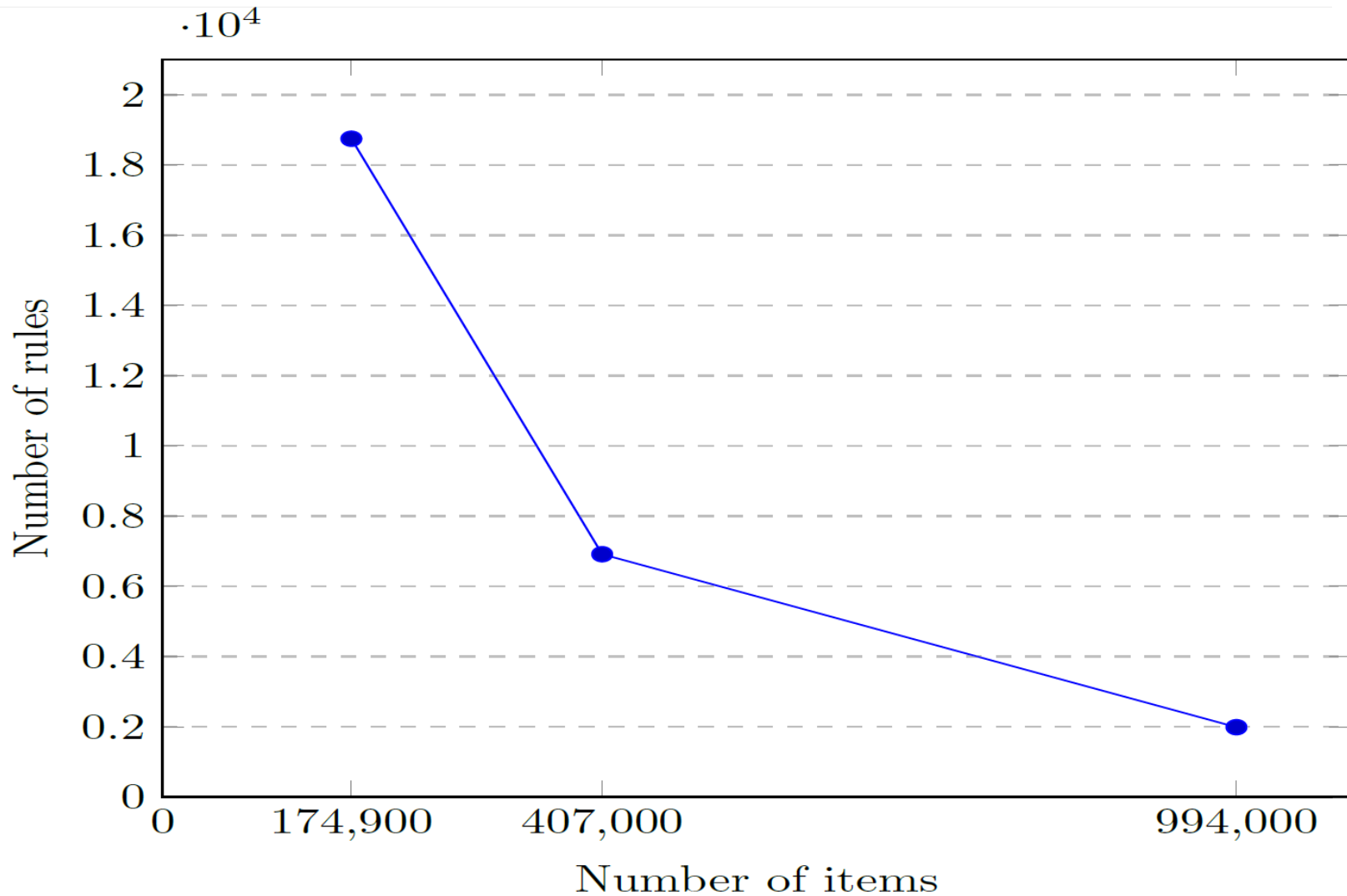| Parameter | Step | Value |
|---|---|---|
| Alpha Threshold | Resolution | 90 |
| Beta Threshold | | 80 |
| Term ranking (n) | | 5 |
| p | | 0.25 |
| Term Threshold ($\gamma$) | Realization | 2 |

Number of **classifications:** $Item \sqsubseteq \forall isClassified.Label$
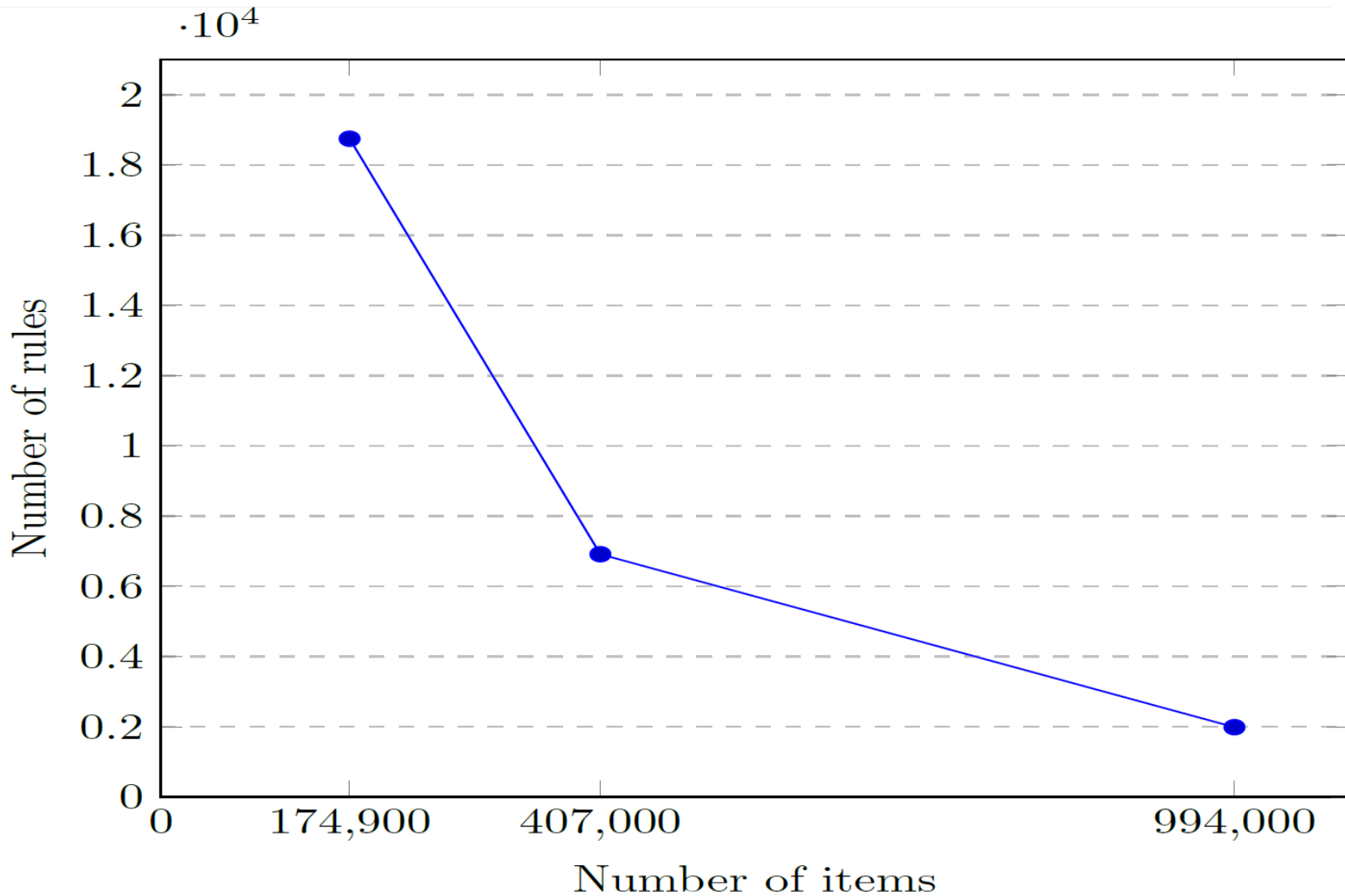
## Number of **learned rules** (Alpha + Beta)

# Number of **learned rules** (Alpha + Beta)  $\alpha$ = 90  $\beta$ = 80

## Context

- Global problem
- The Semantic HMC

## Specific Problem

- Proposed Solution

## Implementation

- Setup
- Results

## Conclusion and future work

# Conclusion

- A new unsupervised process to automatically classify items

  - A highly scalable rule learning method based on statistical and lexical approaches
  - A novel method to classify items using a web reasoner

- The process prototype was successfully implemented in a scalable and distributed platform to process Big Data

- Preliminary results show that the items are classified automatically by the reasonner

- Quality Evaluation of the process: comparison with state-of-the art in classification

- Predictive performance evaluation based on cross-validation with large dataset

- Optimization of the process by exhaustive analysis of parameters' impact

- Application to classification of news articles on the web

# An unsupervised classification process for large datasets using web reasoning

Thank you !

Rafael PEIXOTO, Thomas HASSAN, Christophe CRUZ, Aurelie BERTAUX, Nuno SILVA
thomas.hassan@u-bourgogne.fr

Laboratoire LE2I – UMR CNRS 6306 – Université de Bourgogne