

# Open Corpus User Modeling

Michal Barla  
Ontožúr 18.4.2010





MAVRIANĀ SKALA  
554 m

MAVRIANĀ SKALA	1.00
MAVRIANĀ SKALA	1.00
MAVRIANĀ SKALA	1.00



# Open Corpus User Modeling

It is not a brain surgery...



# Open Corpus **User Modeling**

- collect and analyze data to produce User Model
  - knowledge
  - goals
  - preferences
  - interests
  - experience
  - background
  - ...

# Open Corpus **User Modeling**

- collect and analyze data to produce User Model
  - knowledge
  - goals
  - preferences
  - interests
  - experience
  - background
  - ...

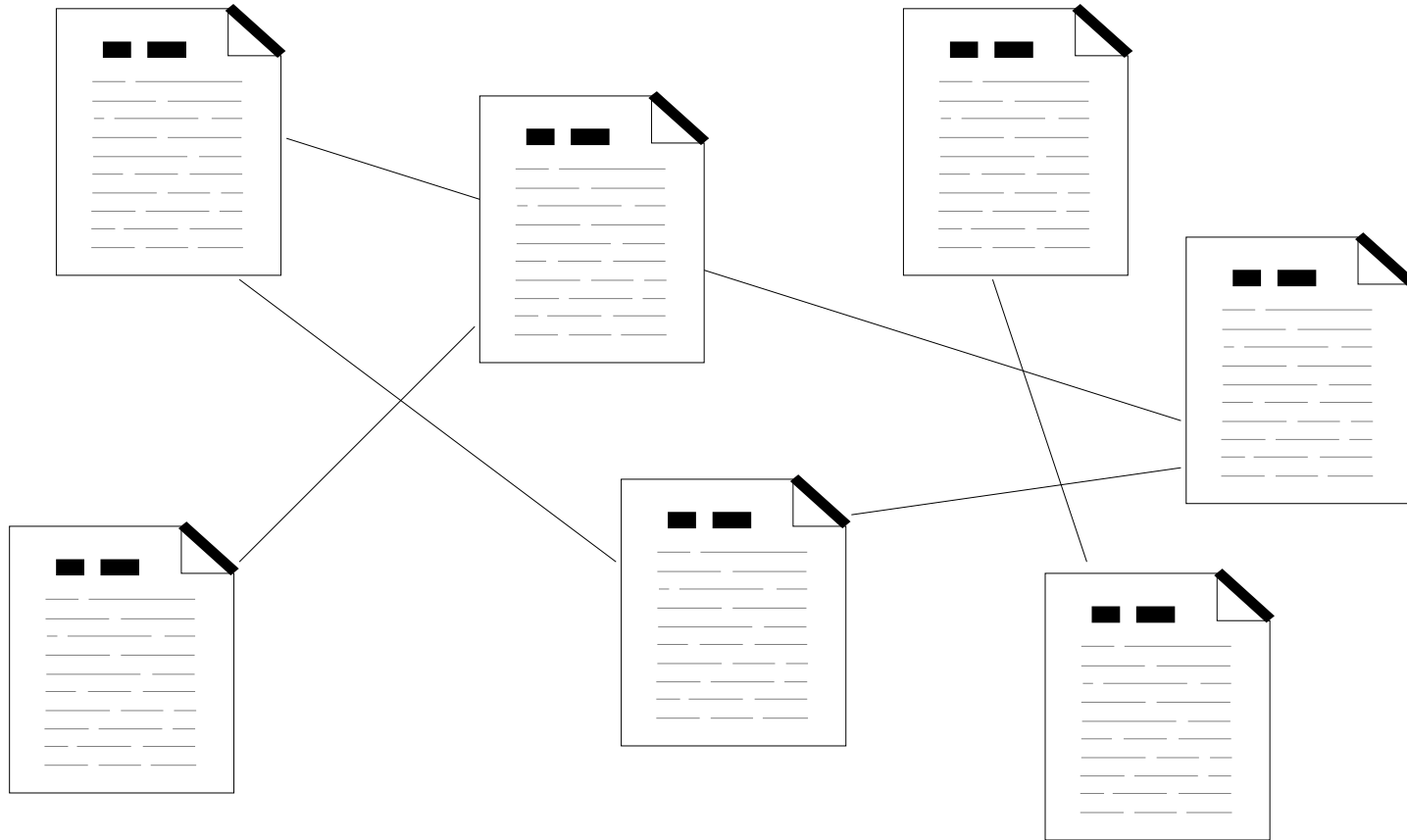
Characteristics  
are related to  
domain elements

# Open Corpus **User Modeling**

- collect and analyze data to produce User Model
  - knowledge
  - goals
  - preferences
  - interests
  - experience
  - background
  - ...

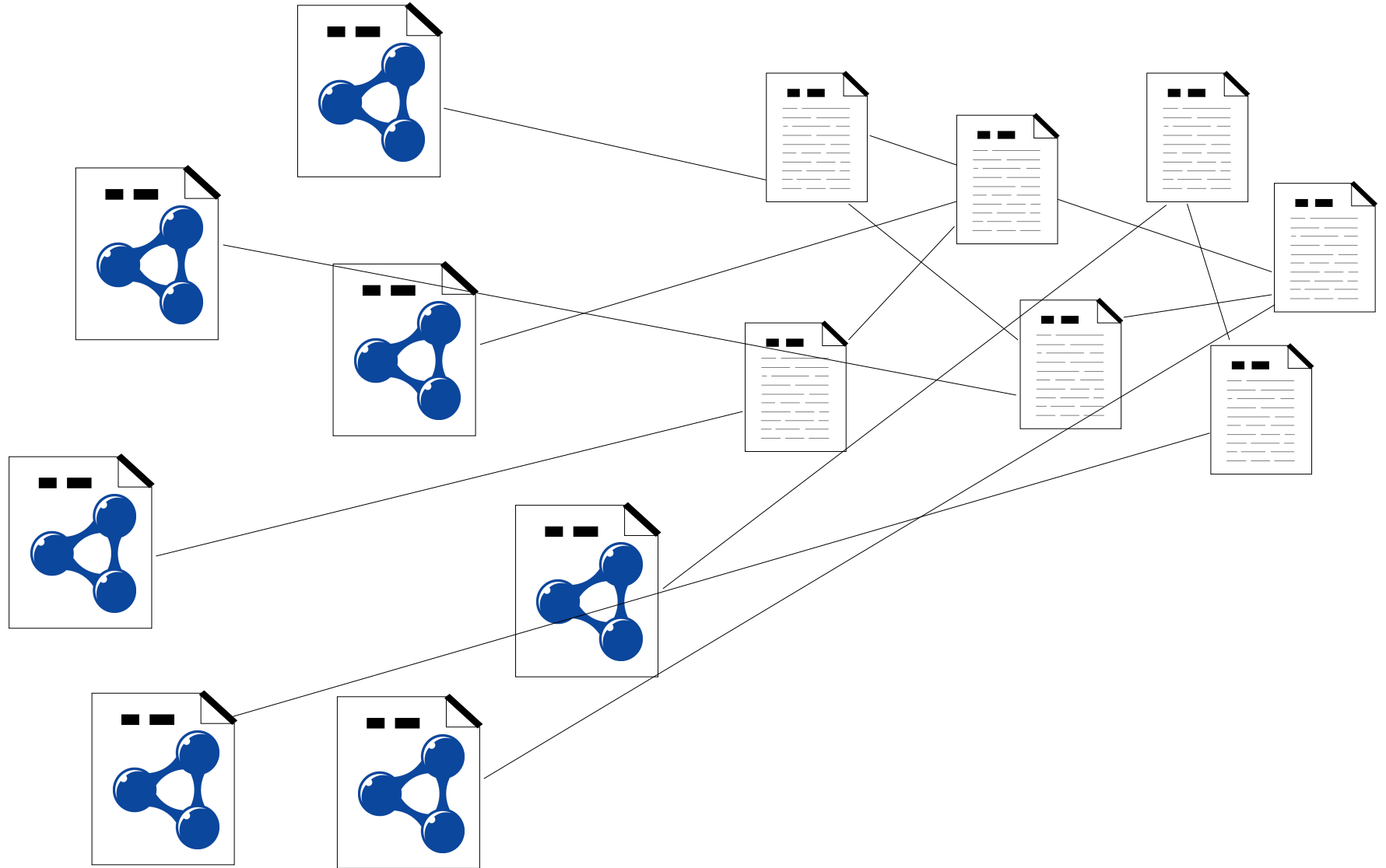
Characteristics  
are related to  
domain elements  
=  
Overlaid UM

# Closed Corpus User Modeling

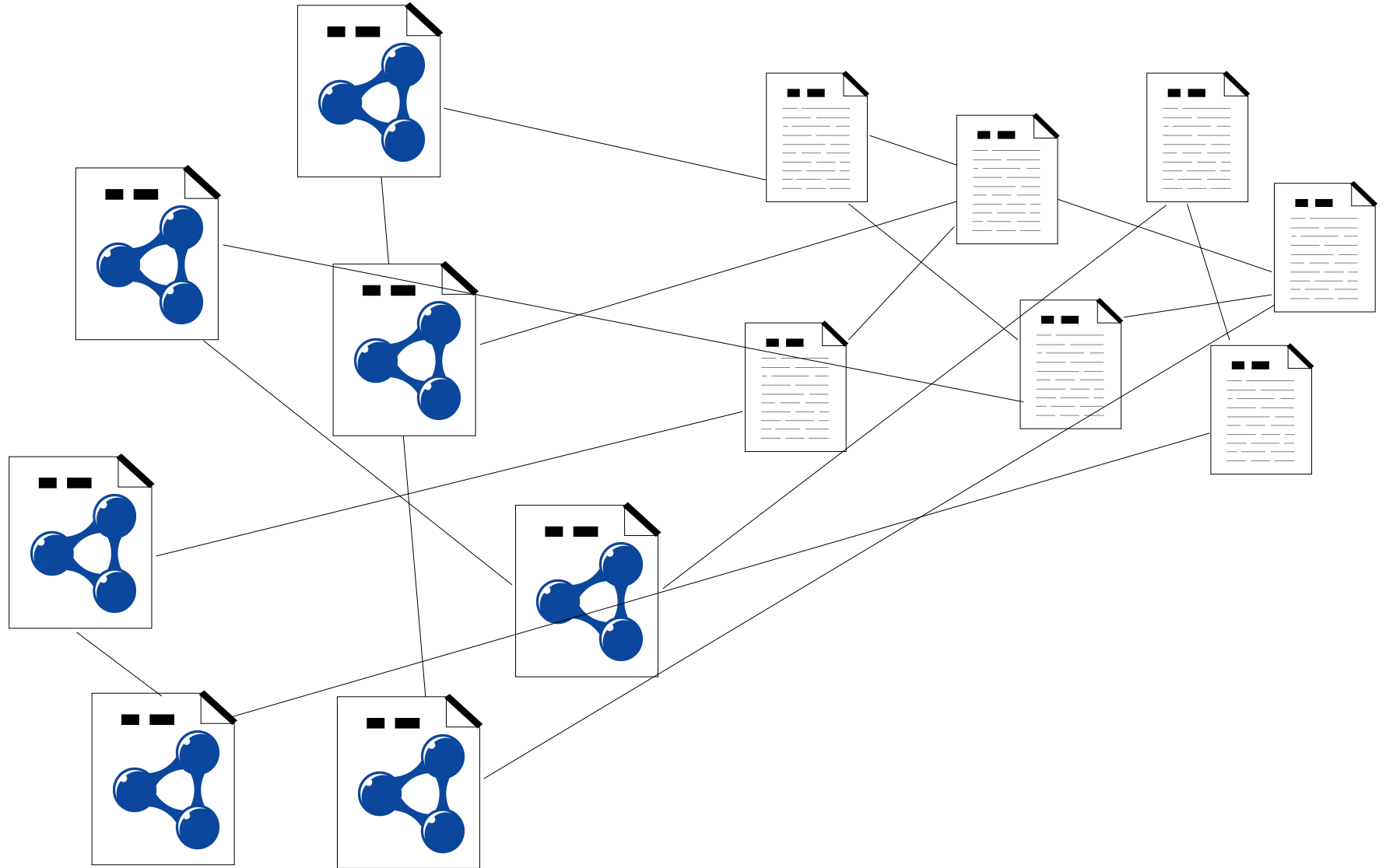




# Open Corpus User Modeling



# Open Corpus User Modeling



# Example 1 – Job offers

Not logged in

Job Offer Portal

User:  Password:

Home Job Offers Factic-TopK Criteria Search Registration User registration About us Help

Start date

- August (2)
- December (1)
- November (7)
- October (36)

Current restrictions

Start date: [All](#) > 2005 (46)  
Duty location: [All](#) > [World](#) > [America](#) > North America (46)  
Offered position: [All](#) > [Professionals](#) > Physical mathematical and engineering science professionals (46)

Duty location

- United States (46)

Sort by: [Name](#) | [Salary](#) | [Organization](#) | [Region](#) Item per page: [10](#) | [15](#) | [25](#) | [50](#) | [100](#)

Offered position

- Computing professionals (46)

Offered management level

- team leader (1)
- worker (1)

Acquisition date

- Last 2 Years (6)
- Last 5 Years (6)

Hours/week

- 30.0 - 40.0 (14)

Travelling involved

#	Name	Salary	Organization	Region	Rate it!
9,38	<a href="#">PeopleSoft Programmer</a>	53.81	Manpower Professional	Durham	☆☆☆☆☆☆
9,02	<a href="#">Programmer; VC++ Makefile generation</a>	30.0		Phoenix	☆☆☆☆☆☆
8,26	<a href="#">Genesis10, Lake Oswego, OR US</a>	35.0	Genesis10	Oregon	☆☆☆☆☆☆
8,05	<a href="#">Programmer Analyst Needed In Chicago!</a>	45000.0		Chicago	☆☆☆☆☆☆
8,00		60000.0		Erie	☆☆☆☆☆☆
7,95	<a href="#">Programmer</a>	29.7		Fort Lauerdale	☆☆☆☆☆☆
7,59	<a href="#">SAS Programmer</a>	35.0		Philadelphia	☆☆☆☆☆☆
7,43	<a href="#">Programmer</a>	25.0	Manpower Professional	Panama City	☆☆☆☆☆☆
7,36	<a href="#">Sr. Oracle P/A with Accounting Bkgmd</a>	75000.0	Diversified Technical Solutions, Inc.	Summit	☆☆☆☆☆☆
7,02	<a href="#">Software Programmer; VC++ Makefile generation</a>	30.0	Manpower Professional	Phoenix	☆☆☆☆☆☆
7,00	<a href="#">Database Programmer- VB6/.net/SQL</a>	23.0	Manpower Professional	Seymour	☆☆☆☆☆☆
6,96	<a href="#">VB.Net Programmer Analyst</a>	55000.0		Columbia	☆☆☆☆☆☆
6,82	<a href="#">HMS Associates Of Tri-State Inc.</a>	55.0	HMS Associates Of Tri-State Inc.	New York	☆☆☆☆☆☆
6,74	<a href="#">ORACLE APPLICATION PROGRAMMER</a>		ATMI	Danbury	☆☆☆☆☆☆

# Example 1 – Job offers

- Each job offer has
  - duty location
  - salary
  - position
  - requirements
- Instead of modeling user's attitudes towards particular job offer, we capture user preferences of different types of attributes



# Example 2 – ALEF

Ste prihlásený ako používateľ Michal Barla (barla, 3653). [Odhlásiť](#)

## Odporúčame pozrieť:

- [Príklad nový](#)
- [Príklad sucet](#)
- [Príklad prvych\\_k](#)
- [Príklad parne](#)
- [Príklad bez\\_tri](#)

[Texty](#) [Otázky](#) [Cvičenia](#)

[\=](#)  
[\==](#)  
[<](#)  
[=<](#)  
[>](#)  
[>=](#)  
[atom](#)  
[číslo](#)  
[filter \(p.t.\)](#)  
[hľadanie \(p.t.\)](#)  
[is](#)  
[member](#)

## [Cvičenie] Príklad vloz

### Zadanie:

Naprogramujte predikát `vloz(+Zoz1, +Zoz2, ?Vysledok)`, ktorý sa splní, ak zoznam `Vysledok` zodpovedá zlúčeniu dvoch usporiadaných zoznamov čísiel `Zoz1` a `Zoz2`, pričom prvky sa nesmú opakovať.

```
?- vloz([1,2,5,8], [1,3,5,7,8], V).  
V = [1,2,3,5,7,8] ->;  
no
```

[Poznám riešenie](#)

[Nepoznám riešenie](#)

[Predchádzajúci príklad / otázka](#)

[Ďalší príklad / otázka](#)

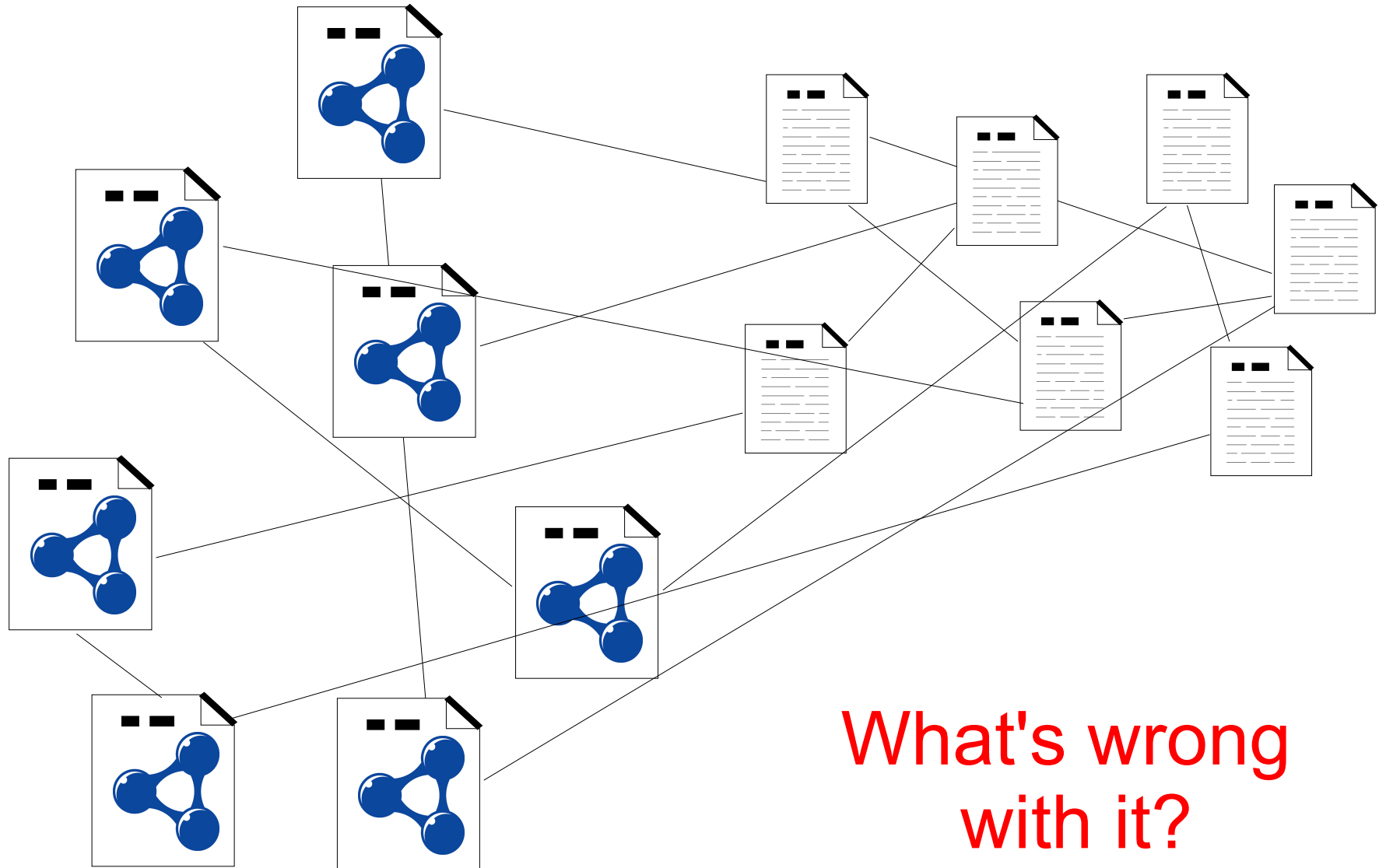
# Example 2 - ALEF

- User is presented with learning objects (LO)
  - explanations
  - questions
  - exercises
- Learning objects are mapped to concepts
- Various concept-to-concept relationships
  - prerequisite
  - parent – child
- We model user knowledge of particular concepts

# Open Corpus User Modeling

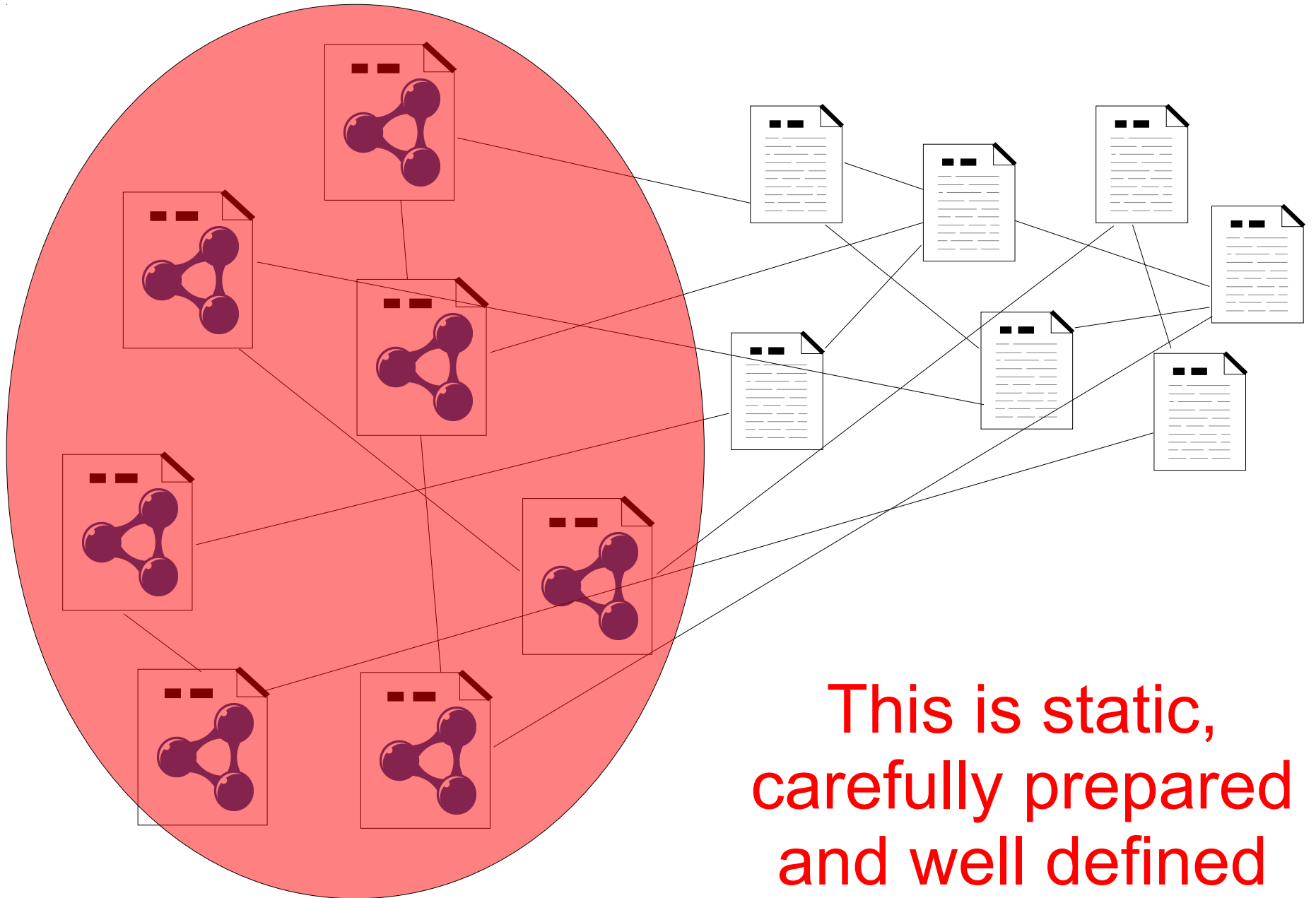
- Instead of overlaying the domain items, we put the user model layer on the top of domain items conceptualization
- Re-use of user model across similar domains
  - Loops in C# and loops in Java
- Support changes of underlying domain items without losing user-related information
  - all we need to do is to map new content to our conceptualization

# Open Corpus User Modeling





# Open Corpus User Modeling



**This is static,  
carefully prepared  
and well defined**

# Static Conceptualization

- Limited to a particular domain
- Personalization is limited to isolated applications
  - Or isolated groups of applications sharing one UM
- Small islands within the whole web ocean

# Adapt. & Person. of “the Web”

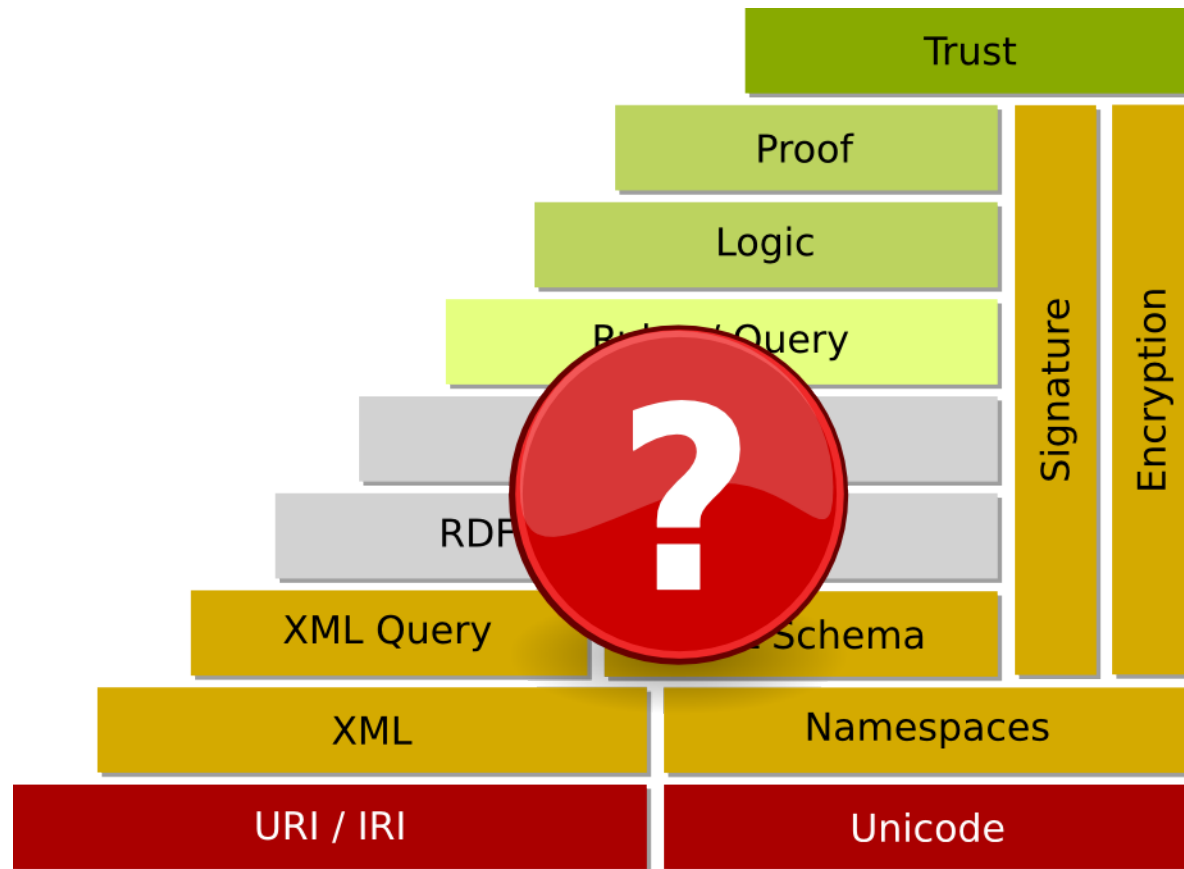
- Pure social-based approach
  - does not require any domain representation
  - System does not need to know the content, users will “tell” it whether it is worthy or not
  - Partitioning users into communities helps to provide more personalized recommendations, annotations etc. (clickstream analysis)

# Adapt. & Person. of “the Web”

- Open corpus domain representation
  - dynamic and open conceptualization
  - Capability to retrieve and process metadata to any document (page) being viewed
- How can this be achieved?



# Can we take advantage of semantic data?



# Problems of the Semantic Web

- Existing semantic systems restricted to a limited set of domains
  - a priori defined domain specific ontologies
  - no links to other ontologies
- The overall Semantic Web does not adequately cover specific terminology
- Many online ontologies have a weak internal structure
  - few online ontologies contain synonyms or non-taxonomic relations,

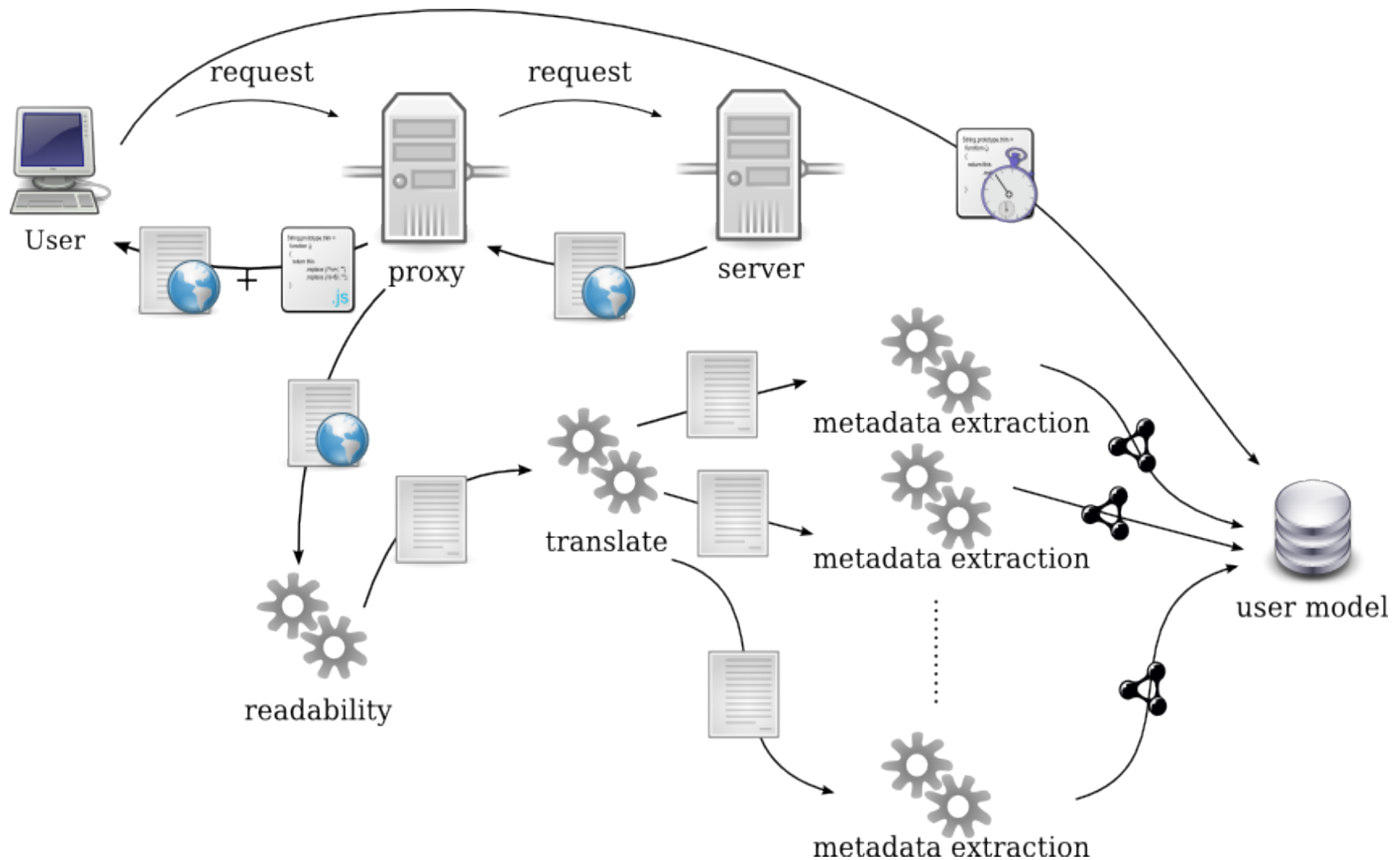
# Waiting for Semantic Web killer app



# Is there any other, reachable semantics?

- People got used to keywords
  - for searching
- People got used to tags and tagging
  - for future retrieval of information
  - but also for first-time retrieval in folksonomies
- People find enough semantics in keywords and tags
- We should give it a try as well...

# Open Corpus **keyword-based** UM



# Metadata extraction

- combining NLP, NER, Linked data and various text processing services
- JATR library
- OpenCalais
- TagTheNet
- AlchemyAPI



# Do you want your own tag-cloud? :)

- Currently 25 distinct IDs
  - max 22 real users
  - top user had 41.693 requests (at the time of preparing this presentation)
  - average was 7693.8
- [peweproxy.fiit.stuba.sk](http://peweproxy.fiit.stuba.sk)
  - query expansion is in alpha testing and looks promising



# Further processing of keywords

- Elimination of stopwords
- Identification of synonyms
- Clusterization
  - semantic-relatedness (Wordnet, Folksonomies, ...)
  - this gives us different profiles of a user
    - researcher, photographer, programmer, hiker, ...

# What we can do with it?

- Detecting virtual communities
  - “smarter” social-based recommendations
- search query disambiguation
- optimization of search results list
- ...

# Conclusions

- We move from basic user models towards two (and more) layers of user models
  - domain items and their conceptualization
- Conceptualization does not need to be closed either
  - keywords could be good enough for performing cool stuff on the top of the “wild wild Web”
- Requires analysis of web traffic
  - proxy server
  - client-side agent as a part of huge multi-agent system