

User Modeling and Personalization

4: User Modeling Frameworks and Interoperability

Eelco Herder

L3S Research Center / Leibniz University of Hanover
Hannover, Germany

12 May 2014

Outline I

Generic User Modeling Systems

- Advantages, Characteristics and Services

- Decentralized Approaches

- In more detail: GUMF

The Semantic Web and the Web of Data

- RDF

- Ontologies

- Some popular ontologies

- Linked Data

User Model Interoperability

- Lingua Franca Approaches

- Conversion Approaches

- Aggregation of User Models

User Modeling Systems

Generic User Modeling Systems

A user modeling system is defined as a 'generic user modeling system' (GUMS) if it is independent from the architecture and from the user model of a specific user-adaptive application.

Generic user modeling systems serve as a separate user modeling component in an application system at runtime.

They may be part of a local area network or a wide area network and serve more than one application instance at a time.

Advantages of GUMS

- ▶ All information about the user is maintained in a repository with clearly defined points of access
- ▶ User information is at the disposal of more than one application at a time
- ▶ User information acquired by one application can be employed by other applications, and vice versa
- ▶ Information about users is stored in a non-redundant manner

(continued on next slide)

- ▶ It is easier to maintain consistency and coherence of information gathered by different applications
- ▶ Information about user groups (e.g. stereotypes) can be maintained with low redundancy
- ▶ Methods and tools for system security, identification, authentication, access control and encryption can be applied and maintained
- ▶ Privacy, transparency

Requirements/Characteristics of GUMS

Several characteristics of generic user modeling systems have been regarded as very important over the years:

- ▶ **Generality:** the system should be usable in as many domains as possible
- ▶ **Expressiveness:** different kinds of user characteristics should be stored in various suitable formats
- ▶ **Inferential capabilities:** GUMS are expected to support complex assumptions and complex reasoning, to be useful in a wide range of domains
- ▶ **Quick Adaptation:** Provide reasonable (or default) assumptions even if there is only a limited amount of user data available

- ▶ **Extensibility:** Interfaces that allow for the (possibly bi-directional) exchange of user information
- ▶ **Import functionality:** Existing data (for example LDAP) should be accessible to the user modeling server
- ▶ **Management of distributed information:** Integration of several sources of user information, such as user profiles, purchase records, ...
- ▶ **Support for open standards:** such as LDAP, FOAF, OpenID

- ▶ **Load balancing:** response delays or denials of requests should only occur in emergency situations
- ▶ **Failover strategies:** fallback mechanisms in case of a breakdown
- ▶ **Transactional consistency:** avoid inconsistencies due to parallel read/write on the user model
- ▶ **Privacy support:** adhere to privacy policies, provide privacy-enhancing services

Services of GUMS

The following services are frequently found in GUMS:

User Data Acquisition

- ▶ The recording of users' behavior, particularly interaction with the system

Knowledge Inference

- ▶ The formation of assumptions based on the interaction history
- ▶ The generalization of the interaction histories of many users into stereotypes
- ▶ The drawing of additional assumptions about the current user based on initial ones

User Model Representation

- ▶ The representation of assumptions about one or more types of user characteristics
- ▶ The representation of relevant common characteristics in specific user subgroups (i.e. stereotypes)
- ▶ The classification of users to a subgroup or stereotype

User Model Management

- ▶ Consistency maintenance in the user model
- ▶ The provision of current assumptions about the user, as well as justifications for these assumptions
- ▶ The evaluation of the entries in the current user model and the comparison with given standards

Decentralized Approaches

Centralized GUMS have several disadvantages:

- ▶ GUMS may impose a restrictive centralized model, or insufficient interfaces for the exchange of user data
- ▶ Reliability and load balancing may be improved by introducing mirrors, but this leads to problems with synchronization and coordination
- ▶ Who is responsible for the protection of the user data (the client adaptive systems or the GUMS administrators)?

Decentralized User Modeling Systems

In decentralized user modeling approaches, each adaptive system maintains its own user model, as needed for its own purposes. The systems communicate directly in a peer-to-peer manner to exchange user profiles and/or user data.

Decentralized user modeling solutions typically offer

- ▶ functionality for the mapping and integration of different types of models
- ▶ solutions for the discovery of communication partners and services

Mixed Approaches

Mixed centralized and decentralized approaches exist as well:

- ▶ user models are stored locally (decentral), but adhering to a common format
- ▶ user models are stored in a decentralized way, but mediation is done through a centralized point of access

Centralized and decentralized approaches to GUMS

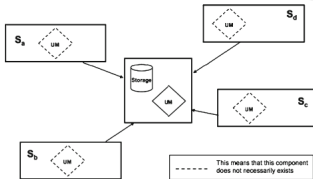


Fig. 1 Centralized approach to UM interoperability

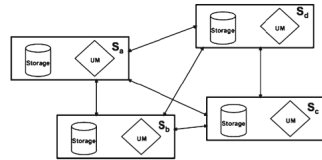


Fig. 2 Decentralized approach to UM interoperability

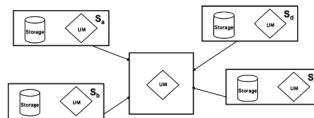


Fig. 3 Mixed approach to UM interoperability

Figure: Centralized, decentralized and mixed UM approaches

In more detail: GUMF

The Grapple User Modeling Framework is a repository for data about users

- ▶ Client applications can query for data and store data
- ▶ GUMF data adheres to the notion of Grapple statements.
- ▶ A Grapple statement is a statement about a user enriched with some metadata
- ▶ Data is logically organized in dataspaces
- ▶ Dataspaces can be enhanced with plug-ins that can provide reasoning functionality or simply gather additional data

GUMF Architecture

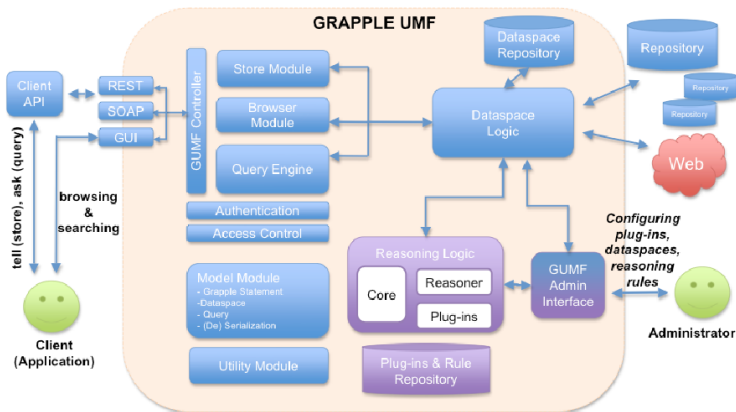


Figure: Conceptual architecture of GUMF

Enriching user profiles

Tim Berners-Lee

The Semantic Web is an extension of the current web in which information is given well-defined meaning, better enabling computers and people to work in cooperation.



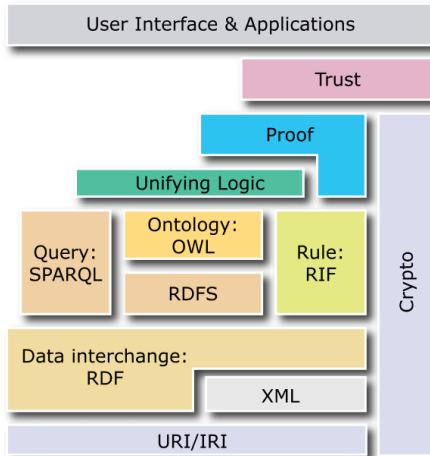
The Semantic Web

The Semantic Web provides a common framework that allows data to be shared and reused accross application, enerpise and community boundaries.

It is a collaborative effort led by W3C with participation from a large number of researchers and industrial partners.

It is based on the Resource Description Framework (RDF), which integrates a variety of applications using XML for syntax and URIs for naming.

Semantic Web Layers



RDF

RDF is a model for representing information about resources in the WWW.

In particular: metadata about Web resources (e.g., title, author, version, publication date, ...)

RDF can also be used as a model for representing information about things that can be identified in the WWW.

In particular: metadata about items (e.g., prices, specifications, availability information, etc.)

In RDF, things are identified using Web Identifiers (Uniform Resource Identifications, URIs)

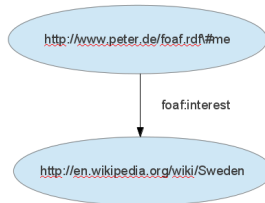
Things are described using properties and property values

The core building blocks of RDF are *triples*: subject - predicate - object.

‘Peter is interested in Sweden’

subject: <http://www.peter.de/foaf.rdf#me>;
predicate: foaf:interest;
object: <http://en.wikipedia.org/wiki/Sweden>;

This statement can be represented as a graph.



RDF Schema and ontologies

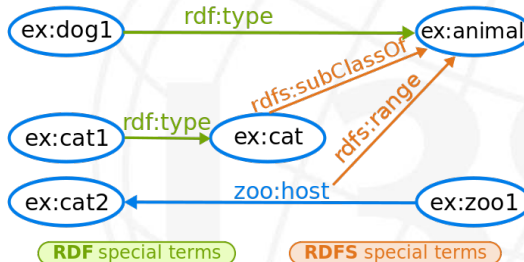
RDF Schema (RDFS) is a set of classes with certain properties using the RDF extensible knowledge representation language, providing basic elements for the description of ontologies.

An **ontology** formally represents knowledge as a set of concepts within a domain, and the relationships between pairs of concepts. It can be used to model a domain and support reasoning about entities.

For example, if we want to represent the following statements:

- ▶ Dog1 is an animal
- ▶ Cat1 is a cat
- ▶ Cats are animals
- ▶ Zoos host animals
- ▶ Zoo1 hosts the Cat2

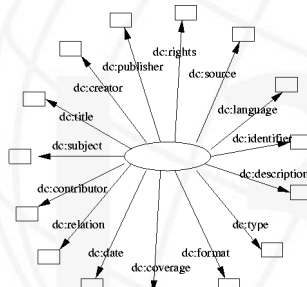
This would result in the following RDF graph:



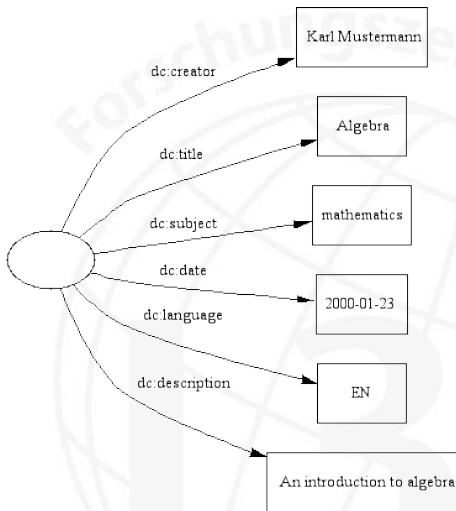
Dublin Core

The Dublin Core metadata terms are a set of vocabulary terms which can be used to describe resources for the purposes of discovery.

The terms can be used to describe a full range of web resources (video, images, web pages, etc.), physical resources such as books and objects like artworks.



An example Dublin Core description of a book.



```
<?xml version="1.0"?>
<rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax
        xmlns:dc="http://purl.org/dc/elements/1.1/">
  <rdf:Description>
    <dc:creator>Karl Mustermann</dc:creator>
    <dc:title>Algebra</dc:title>
    <dc:subject>mathematics</dc:subject>
    <dc:date>2000-01-23</dc:date>
    <dc:language>EN</dc:language>
    <dc:description>An introduction to algebra</dc:descript
  </rdf:Description>
</rdf:RDF>
```

Friend Of A Friend (FOAF)

FOAF is a language to describe people and resources in the Web.

FOAF Basics	Personal Info	Online Accounts / IM	Projects and Groups	Documents and Images
<ul style="list-style-type: none"> • Agent • Person • name • nick • title • homepage • mbox • mbox_sha1sum • img • depiction (depicts) • surname • family_name • givenname • firstName 	<ul style="list-style-type: none"> • weblog • knows • interest • currentProject • pastProject • plan • based_near • workplaceHomepage • workInfoHomepage • schoolHomepage • topic_interest • publications • geekcode • myersBriggs • dnaChecksum 	<ul style="list-style-type: none"> • OnlineAccount • OnlineChatAccount • OnlineEcommerceAccount • OnlineGamingAccount • holdsAccount • accountServiceHomepage • accountName • icqChatID • msnChatID • aimChatID • jabberID • yahooChatID 	<ul style="list-style-type: none"> • Project • Organization • Group • member • membershipClass • fundedBy • theme 	<ul style="list-style-type: none"> • Document • Image • PersonalProfileDocument • topic (page) • primaryTopic • tipjar • sha1 • made (maker) • thumbnail • logo

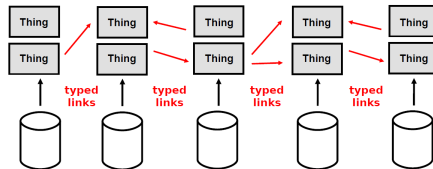
Linked data

The Web enables us to link related documents. Similarly it enables us to link related data.

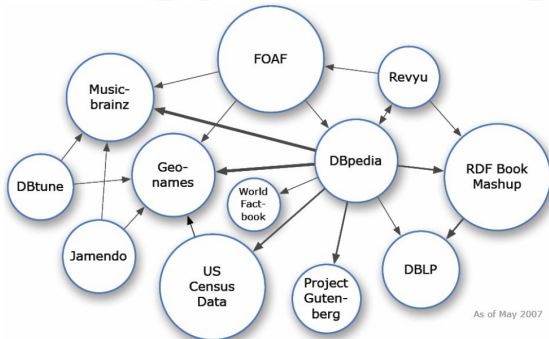
Linked Data refers to a set of best practices for publishing and connecting structured data on the Web.

Key technologies that support Linked Data are

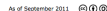
- ▶ **URIs** - a generic means to identify entities or concepts in the world
- ▶ **HTTP** - a simple yet universal mechanism for retrieving resources, or descriptions of resources
- ▶ **RDF** - a generic graph-based data model with which to structure and link data that describes things in the world



The Linking Open Data (LOD) project takes existing open data sets, makes them available on the Web in RDF and interlink them with other data sets.



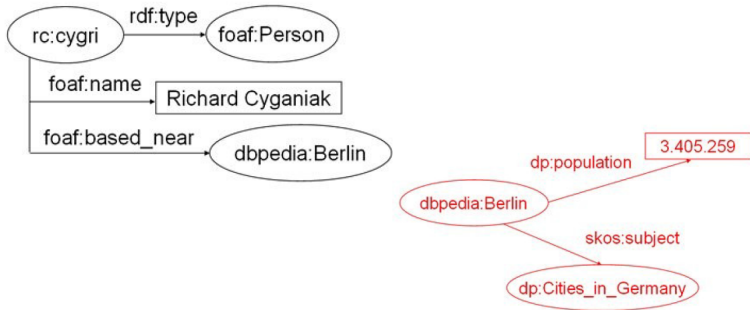
2007



Eelco Herder | User Modeling and Personalization 4: User Modeling Frameworks and Interoperability | 32/55

Making use of linked data principles and the “LOD Cloud”, several pieces of knowledge can be connected.

For example, making use of DBpedia, we can infer that Richard Cyganiak lives in a city with a population of 3.405.259.



User Model Interoperability

Interoperability

The ability to cooperate and exchange data despite differences in languages, interface and execution platform

The main advantages of interoperable user models are:

- ▶ acquiring more data and more accurate data about users
- ▶ acquiring functionalities (both user model and user modeling functionalities) that systems do not themselves implement

Cross-application user model interoperability is a strategy to cope with the *cold start problem*, which refers to the difficulty for applications to provide suitable adaptations for new users.

Interoperability may also relieve users from the pain of training new systems or wasting time filling in their profile for every new application.

In general, it is believed that interoperable user models lead to more user data and more accurate models.

Dimensions of Interoperability

In order to successfully exchange and (re-)use user profile data, requirements on several levels need to be met:

- ▶ *Structural level*: Protocols for communication (i.e. interfaces)
- ▶ *Syntactic level*: Languages for the exchange of user model data
- ▶ *Semantic level*: Agreement on definitions what certain concepts (such as knowledge), terms and expressions, exactly mean

Protocols for communication are abundant:

- ▶ Remote calls (remote procedure calls, Java remote service invocation, other client-server approaches)
- ▶ Web-based protocols (HTTP POST/GET, WSDL/SOAP, XML, JSON, OAI, RESTFUL interfaces)

Syntactic and semantic interoperability is much harder, as we will see in the next part of this lecture.

Two Basic Approaches

In essence, there are two ways to ensure interoperability between two (or more) adaptive systems and their user models:

- ▶ **Lingua franca:** an agreement between all parties on a common representation and semantics
This is the philosophy underlying the generic (centralized) user model server approach
- ▶ **Conversion** of user model information between the different systems.

Conversion allows for flexible and extensible user models, at the price of possible loss of information in the conversion process:

- ▶ models may be simply incompatible (there is no suitable mapping)
- ▶ mappings may be incomplete (information required in one model is not available in the other)
- ▶ the observations in the different systems may lead to contradictions

Lingua Franca Approaches

Generic User Modeling Ontology

GUMO is developed at the DFKI German Institute for Artificial Intelligence and is divided in four parts:

- ▶ *GUMO-Basic*: models user dimensions like personality traits and user characteristics
- ▶ *GUMO-Context*: models the world directly around the user, including the events and the environmental context
- ▶ *GUMO-Domain*: includes a general interest (domain overlay) model
- ▶ *GUMO-Extended*: ranges, rating dimensions and predicates that address other attributes

By adding new parts to GUMO, other ontologies (for example elearning-related) can be integrated.

GUMO

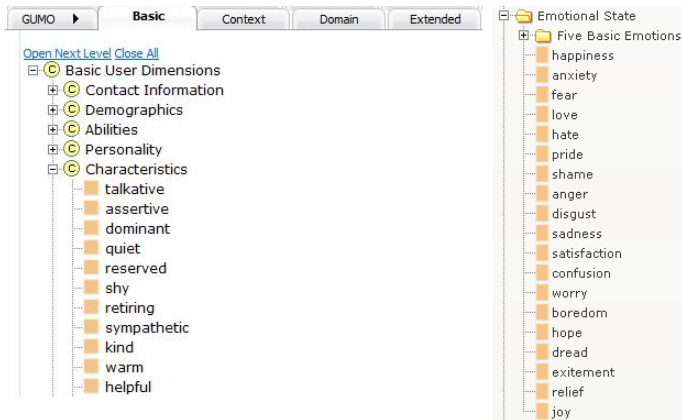


Figure: Screenshot of GUMOs user interface

Conversion Approaches

We illustrate conversion approaches by means of a Grapple Derivation Rules that convert data from one format into another.

This simple rule maps test scores of the elearning system CLIX to a more generic external knowledge format by simply multiplying the score by 10:

```
<gdr:rule
  xmlns:gdr="http://www.grapple-project.org/grapple-derivation-rule/"
  xmlns:gc="http://www.grapple-project.org/grapple-core/"
  xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"

  id="1"
  name="Quiz Level to External Knowledge (CLIX)"
  description="Maps the quiz/test result of a specific CLIX score to external knowledge"
  creator="http://pcwin530.win.tue.nl:8080/grapple-umf/client/1">
    <gdr:premise dataspace="1">
      <gc:subject>?user</gc:subject>
      <gc:predicate rdf:resource="http://www.grapple-project.org/ims-lip/completedTest"/>
      <gc:object>solar-system-quiz-id</gc:object>
      <gc:level>?level</gc:level>
    </gdr:premise>
    <gdr:consequent dataspace="1">
      <gc:subject>?user</gc:subject>
      <gc:predicate rdf:resource="http://gale.tue.nl/predicate/knowledge"/>
      <gc:object>gale://gale.tue.nl/cam/DavidTestCAM/SolarSystem</gc:object>
      <gc:level>op:multiply(?level,10)</gc:level>
    </gdr:consequent>
  </gdr:rule>
```

Aggregation of User Models

Users leave different types of profile traces on the Social Web.

People fill out their profile attributes, such as name, affiliations, etcetera.

Social tagging systems capture tagging activities of the users to create *tag-based profiles*.

Form-based profile

The form-based profile of a user u is a set of attribute-value pairs.

$$UM(u) = \{\{a, v\} | a \in A_{UM} \text{ and } v \text{ is in the range of } a\}$$

A_{UM} defines the vocabulary of attributes that can be applied to describe characteristics of the user u . The value v associated with an attribute a must be in the range of a .

Traditional attributes might be name or email address:

$$\blacktriangleright UM(u_1) = \{(name, 'Bob'), (email, bob@mail.com)\}$$

Tag-based profile

The *tag-based profile* of a user u is a set of weighted tags where the weight of a tag t is computed by a certain strategy w with respect to the given user u .

$$P(u) = \{\{t, w(u, t)\} | t \in T_{source}, u \in U\}$$

$w(u, t)$ is the weight that is associated with tag t for a given user u . T_{source} is the source set of tags from which tags are extracted for the tag-based profile $P(u)$.

For example, $P(u_1) = \{(\text{research}, 0.65), (\text{semantic web}, 0.2), \text{jazz}, 0.15)\}$

Aggregation of form-based profiles

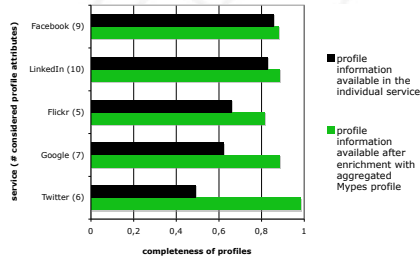
form-based profile attributes	Face- book	LinkedIn	Twitter	Blog- Spot	Flickr	Delicious	Stumble Upon	Last.fm	Google
nickname	x	x	x	x	x	x	x	x	x
first name	x	x							
last name	x	x							
full name	x	x	x		x				x
profile photo	x		x		x				x
about		x							x
email (hash)	x				x				
homepage	x	x	x						x
blog/feed			x	x	x	x	x	x	
location		x	x		x				x
locale	x								
interests		x							
education		x							
affiliations	x	x							
industry		x							

Table: Form-based profile attributes available in Social Web Platforms

Completeness of user profiles

Users often do not fill out their profiles completely. For example, Twitter only asks 6 attributes, but these profiles are only completed up to 49%.

Would it be possible to complete user profiles by aggregating data from different sources?



Fabian Abel, Eelco Herder, Geert-Jan Houben, Nicola Henze, Daniel Krause. Cross-system User Modeling and Personalization on the Social Web. UMUI 23 (2-3), 2013, pp 169-209

Aggregation of tag-based profiles

How useful is a tag-based profile? Intuitively, if the profile contains just one tag, this profile is not very useful. The profile becomes more varied if other tags are included - in particular if there is not one tag that appear 8 out of 10 times in the profile.

One way to measure the usefulness of a tag-based profile is the *entropy* of the collection.

Entropy

Entropy is a measure of disorder, or more precisely unpredictability.

For example, a series of coin tosses with a fair coin has maximum entropy, since there is no way to predict what will come next.

A string of coin tosses with a two-headed coin has zero entropy, since the coin will always come up heads.

Most collections of data in the real world lie somewhere in between.

The entropy of a tag-based profile T , which contains of a set of tags t , is computed as follows:

$$\text{entropy}(T) = \sum_{t \in T} p(t) \cdot \text{self-information}(t) \quad (1)$$

In Equation 1, $p(t)$ denotes the probability that the tag t was used by the corresponding user. Self-information is the logarithm of $p(t)$ multiplied by -1 :

$$\text{self-information}(t) = -\log_2(p(t)) \quad (2)$$

Using base 2 for the computation of the logarithm allows for measuring self-information as well as entropy in bits.

For modeling the probability $p(t)$ that a tag t appears in a given user profile, we apply the individual usage frequencies of the tags, i.e. for a specific user u the usage frequency of tag t is the fraction of u 's tag assignments where u referred to t .

Overlap

Aggregation of tag-based profiles is meant to reveal more information about the users than the profiles available in some specific service.

If you combine two more or less identical profiles, the added value is rather low. If you combine two completely different profiles (e.g. your favorite music with your research interests), the entropy will increase, but the aggregated profile is probably less useful.

Therefore, there needs to be at least *some* overlap between the profiles that you combine.

For each user u and each pair of services A and B , we compute the overlap as specified in Definition 3.

$$\text{overlap}(u_A, u_B) = \frac{1}{2} \cdot \left(\frac{|T_{u,A} \cap T_{u,B}|}{|T_{u,A}|} + \frac{|T_{u,A} \cap T_{u,B}|}{|T_{u,B}|} \right) \quad (3)$$

$T_{u,A}$ and $T_{u,B}$ denote the set of (distinct!) tags that occur in the tag-based profile of user u in service A and B respectively. Hence, $|T_{u,A} \cap T_{u,B}|$ is the number of (distinct) tags that occur in both profiles, u_A and u_B .

Example overlaps

