

Relaxing RDF queries based on user and domain preferences

Peter Dolog · Heiner Stuckenschmidt ·
Holger Wache · Jörg Diederich

Received: 17 April 2008 / Accepted: 22 May 2008 /
Published online: 30 June 2008
© Springer Science + Business Media, LLC 2008

Abstract Research in cooperative query answering is triggered by the observation that users are often not able to correctly formulate queries to databases such that they return the intended result. Due to lacking knowledge about the contents and the structure of a database, users will often only be able to provide very broad queries. Existing methods for automatically refining such queries based on user profiles often overshoot the target resulting in queries that do not return any answer. In this article, we investigate methods for automatically relaxing such over-constrained queries based on domain knowledge and user preferences. We describe a framework for information access that combines query refinement and relaxation in order to provide robust, personalized access to heterogeneous resource description framework data as well as an implementation in terms of rewriting rules and explain its application in the context of e-learning systems.

P. Dolog (✉)
Department of Computer Science, Aalborg University,
Selma Lagerlöfs Vej 300, 9220 Aalborg, Denmark
e-mail: dolog@cs.aau.dk

H. Stuckenschmidt
Universität Mannheim - Institut für Informatik,
A5, 6, 68159 Mannheim, Germany
e-mail: heiner@informatik.uni-mannheim.de

H. Wache
School of Business, University of Applied Sciences
Northwestern Switzerland (FHNW),
Riggenbachstrasse 16, CH-4600 Olten, Switzerland
e-mail: holger.wache@fhnw.ch

J. Diederich
Forschungszentrum L3S, Leibniz Universität Hannover,
Appelstraße 9a, 30167 Hannover, Germany
e-mail: diederich@L3S.de

Keywords Query relaxation · User modeling · Preferences · ECA rules

1 Introduction

Cooperative query answering is triggered by the observation that users are often not able to correctly formulate queries to databases that return the intended result. This is even more the case for semantic web systems based on resource description framework (RDF) for the following reasons:

- The underlying data often comes from different sources. The internal structure of these sources is not always known.
- There is no fixed integrated schema. Each source can have its own schema, sources may make partial use of different available schemas.
- Authors of data are not forced to stick to a given schema and often use different conventions to represent the same information.

With the increasing popularity of RDF as a representation language in domains such as medicine (Stuckenschmidt et al. 2004) or e-learning (Dolog et al. 2004, 2008) this problem becomes more pressing. We have to make sure that people will be able to formulate meaningful queries. If this is not the case, we have to find ways to still provide the user with the intended results. Cooperative query processing aims at supporting the user by automatically modifying the query in order to better fit the real intention of the user. The specific problem that we address in this work is the situation where the content of the information as well as the RDF-based metadata do not have the expected format.

As a motivating example, we use the domain of e-learning. In open e-learning environments such as the one described in Dolog et al. (2004, 2008), learning materials from different authors and different sources are made available online through an RDF-based infrastructure. In particular, each learning resource is described by a combination of different metadata providing details about type, format, and content of the material as well as information about required expertise. In order to capture this metadata, different metadata schemas including standards like Dublin Core¹ and LOM² as well as domain ontologies like the ACM topic hierarchy³ are used (Brase 2005). In addition, information about the expertise of the user is needed to decide whether the material meets the user's level of expertise. The result is a rather complex metadata schema that needs to be queried in order to retrieve relevant material. Due to this complexity, the corresponding query is normally created by the system based on user input and known preferences. An example of such a complex query in the SeRQL query language is shown in Fig. 1.⁴

¹<http://dublincore.org/schemas/rdfs/>.

²<http://kmr.nada.kth.se/el/ims/md-lomrdf.html>.

³<http://www.cs.vu.nl/~heiner/public/SW@VU/classification.daml>.

⁴Namespaces have been omitted for the sake of readability.

Fig. 1 Query extended with user preferences

```

SELECT * FROM
  {Resource} subject {Subject}, {Resource} title {Title},
  {Resource} description {Description}, {Resource} language {Language},
  {Resource} requires {} subject {Prerequisite},
  {User} type {user}, {User} hasPerformance {Performance},
  {Performance} learning_competency {Competence}
WHERE
  Subject Like "inferenceengines" and Prerequisite = Competence and
  Language = de and User = user42

```

The query asks for learning resources in German about the subject “inferenceengines” that meet the previous knowledge of the user – in particular, it claims that the subject of required units should be contained in the competence of the current user (user42). When we apply this query to real data, however, it does not return any result despite the fact that there are a number of suitable resources. The problem in this case are irregularities in the way the data is described. In particular, not all the authors of resource metadata stick to the conventions that come with the metadata schema – a situation that we encounter quite often on the semantic web. For the purpose of this article, we will not systematically analyze the different problems we identified in the dataset but restrict ourselves to two typical examples that we will use to illustrate our approach in the remainder of the article.

The subject assigned to a course does not always correctly summarize the content. In our test data set, for example, if the user provides the keyword “inferenceengines” no resources are returned despite the fact that there are resources for instance about inference implementation and tools for inference. The problem here is that in the case of the first resource the term “inferenceengines” only occurs in the title, but not in the subject. In the case of the second resource, the term only occurs in the description and is mentioned neither in the subject nor in the title of the resource. Therefore, a query for resources with the subject “inferenceengines” will return no results.

Learning resources normally specify required expertise in terms of knowledge that the user should have in order to be able to understand the content of the resource. In existing data sets there are at least three different ways in which this requirement is specified. The standard way is to refer to other learning resources that should have been mastered before - this information can be taken from the user profile. There are also cases, however, where required skills are given in terms of links to a topic hierarchy or even in terms of literals containing the required skill as a string. If the user query is now automatically expanded with links to resources mastered by the user, these resources will not be found.

In this article, we present an approach for addressing the problems querying heterogeneous RDF models using preference-based query relaxation. Our work is similar to the work of Gaasterland (Gaasterland et al. 1992a, b) on cooperative query answering in databases. In particular, we present a method for automatically relaxing over-constrained queries based on background knowledge about the domain model and user preferences, extending previous work in two directions:

- We base our work on a general relaxation framework that is designed to capture the specific properties of RDF and RDF query languages.
- We explicitly link the relaxation process to an explicit model of user preferences and domain characteristics that can be used to guide the relaxation process.

In the following, we provide an overview of a general approach for relaxing RDF queries using rewritings that we proposed in previous work. In Section 3, we introduce a generic model for representing domain and user knowledge that provides the basis for guiding the relaxation process which is discussed in detail in Section 5. We also discuss the introduced model in the context of related work in Section 5. Section 4 discusses how to elicit knowledge about a user and a domain. We briefly describe a prototypical implementation of the system in Section 6 and conclude with a discussion of achievements and open problems.

2 Relaxation by rewriting

A possible solution to provide users with meaningful results is to successively relax the constraints imposed in the (extended) query. Different techniques for relaxing queries have been proposed in the database area. Gaasterland et al. (1992a) provide a unifying view on different relaxation techniques in terms of replacing subexpressions in the query.

2.1 The rewriting framework

We proposed in previous work (Dolog et al. 2006) a rule-based query rewriting framework for RDF queries independent of a particular query language. The framework is based on the notion of triple patterns (RDF statements that may contain variables) as the basic element of an RDF query and represents RDF queries in terms of three sets: triple patterns that must be matched (mandatory patterns), triple patterns that may be matched (optional triple patterns), and conditions in terms of constraints on the possible assignment of variables in the query patterns.

Rewritings of such queries are described by transformation rules $Q \xrightarrow{R} Q'$ where Q is the original query and Q' the rewritten query. Rewriting rules consists of three parts:

- A matching pattern represented by an RDF query in the sense of the description above. Normally the matching pattern is a part of the original query Q .
- A replacement pattern also represented by an RDF query in the sense of the description above. The replacement pattern replaces the matched pattern resulting in Q' .
- A set of conditions in terms of special predicates that restrict the applicability of the rule by restricting possible assignments of variables in the matching and the replacement pattern.

A rewriting is performed in the following way: First it is checked if a rule is applicable for a query. The matching pattern needs to match a given query Q in the sense that the mandatory and optional patterns are subsets of the corresponding parts of Q . Furthermore, the predicates in the conditions of the rewriting rule have to be satisfied. If a rule is applicable for a query then the matched patterns are removed from Q and replaced by the corresponding parts of the replacement pattern resulting in Q' .

2.2 Query rewriting model

Our rewriting approach relaxes the over-constrained query based on rules, which has the advantage that we start with the strongest possible query that is supposed to return the “best” answers to satisfy most of the conditions. If the returned result set is either empty or contains unsatisfactory results, the query is modified either by replacing or deleting parts of it, or in other words, relaxed. The relaxation should be a continuous step by step, (semi-)automatic process, to provide a user with the possibility to interrupt further relaxations. Before we investigate concrete relaxation strategies in the context of our example domain, we first give a general definition of the framework for rewriting an RDF query.

Each resource is annotated with an RDF description which can be seen as a set of triples (Hayes 2004). A query over these resources consists of triple patterns and a set of conditions that restrict the possible variable bindings in the patterns. Each triple pattern represents a set of triples. The corresponding abstract definition of a query focuses on the essential features of queries over RDF; several concrete query languages are based on these ideas including SeRQL (Broeskstra and Kampman 2004) which we use in our examples in Fig. 1.

Definition 1 (RDF Query) Let \mathcal{T} be a set of terms, \mathcal{V} a set of variables, \mathcal{RN} a set of relation names, and \mathcal{PN} a set of predicate names. The set of possible triple patterns \mathcal{TR} is defined as $\mathcal{TR} \subseteq 2^{(\mathcal{T} \cup \mathcal{V}) \times (\mathcal{RN} \cup \mathcal{V}) \times (\mathcal{T} \cup \mathcal{V})}$. A ground triple (pattern) is a triple pattern which does not contain any variable. A query Q is defined as $\langle M_Q, O_Q, P_Q \rangle$ with M_Q and $O_Q \in \mathcal{TR}$ and $P_Q \subseteq \mathcal{P}$. M_Q is the set of mandatory patterns (patterns that have to be matched by the result), O_Q is a set of optional patterns (patterns that contribute to the result but do not necessarily have to match the result), and \mathcal{P} is the set of predicates. A predicate has a name in \mathcal{PN} and is defined over \mathcal{T} and \mathcal{V} .

The triple patterns M_Q in a query Q determine ground triples from a database. Informally all substitutions τ are answers to Q which maps the triple patterns M_Q to existing ground triples in the database. A substitution τ is defined as a set of pairs (X_i, T_i) and applied as usual:

Definition 2 (RDF Answers) A substitution τ is a set of pairs (X_i, T_i) with $X_i \in \mathcal{V}$ and $T_i \in \mathcal{T} \cup \mathcal{V}$. $\tau(S)$ is generated from S where each appearance of X_i is replaced by T_i for each $(X_i, T_i) \in \tau$. A substitution τ is valid for a RDF query $Q = \langle M_Q, O_Q, P_Q \rangle$ if

- $M = \tau(M_Q)$ where M is a set of ground triples from the database and
- $\tau(P_Q)$ are satisfied.

A valid substitution τ may be extended to optional patterns O_Q , i.e. $\tau(O_Q)$ may also be equal to some ground triple in database. All valid substitutions constitute answers to the query Q .

Using these abstract definitions, the query in Fig. 1 without language preference on German and prerequisites on competences matched with the user's learner performance would be represented as

$$\begin{aligned}
 M_Q &= \{(Resource, subject, Subject), \\
 &\quad (Resource, title, Title) \\
 &\quad (Resource, description, Description)\} \\
 O_Q &= \{\} \\
 P_Q &= \{like(Subject, "inferenceengines")\}
 \end{aligned}$$

where $Resource, Subject, Title, Description \in \mathcal{V}$, as well as $subject, title, description, "inferenceengines" \in \mathcal{T}$ and $like \in \mathcal{PN}$.

Based on the abstract definition of an RDF query, we can now define the notion of a rewriting rule and rewriting process as such. We define rewriting in terms of rewriting rules that take as input parts of a query, in particular triple patterns and conditions, and replace them by different elements.

In our work, we employ the principle of rewriting rules that are inspired by ECA-rules (event-condition-action rules) (Ceri 1992) for continuous relaxation of user queries. A rewriting rule formally consists of three parts: a *pattern*, a *replacement* and some *conditions*. The pattern corresponds to the event, i.e. in our case an occurrence of particular triple patterns or predicates in a query. The replacement contains the terms which will substitute the matched pattern in a query; the replacement can be seen as the action in the ECA principle. Conditions constrain the rewriting and determine when a particular rule can be fired because the rewriting rule can only be applied if the conditions are satisfied. These conditions can be used to define certain relaxation strategies. In particular, we will see later that conditions can be based on user preferences or background knowledge about the domain.

Definition 3 (Rewriting Rule) A rewriting rule R is a 3-tuple $\langle PA, RE, CN \rangle$ where PA and RE are RDF queries according to Definition 1 and CN is a set of predicates.

Conditions are the same constructs which are already introduced for queries. Conditions consist of predicates which constrain possible results. Patterns and replacements formally have the same structure as queries. They also consist of a set of triples and predicates. But patterns normally do not address complete queries but only a subpart of a query. Using this definition we can specify a rewriting rule that extends the simple query in Fig. 1 with the language preference of the example user `user42`.

$$\begin{aligned}
 PA &= \{(Resource, title, Subject)\}, \emptyset, \emptyset \\
 RE &= \{(Resource, title, Subject), \\
 &\quad (Resource, language, Language)\}, \emptyset, \\
 &\quad \{(Language = X)\} \\
 CN &= \{languagePreference(User, X)\}
 \end{aligned}$$

where *languagePreference* is a predicate which looks in the user profile of a user *User* for the language preference. *User* is a variable which will be bound to the id of the current user sending the query, e.g., `user42`. While this example contained a rule for refining a query, we will see later that we can use the same mechanism for defining relaxations on a query.

In general a rewriting rule is applicable to all queries which contain the pattern as a part. The pattern does not have to cover the whole query. Normally it addresses some triples as well as some predicates in the query. In order to write more generic rewriting rules, the pattern must be instantiated which is done by a substitution.

Definition 4 (Pattern Matching) A pattern *PA* of a rewriting rule *R* is applicable to a query $Q = \langle M_Q, O_Q, P_Q \rangle$ if there are subsets $M'_Q \subseteq M_Q, O'_Q \subseteq O_Q$ and $P'_Q \subseteq P_Q$ and a substitution θ with $\langle M'_Q, O'_Q, P'_Q \rangle = \theta(PA)$.

In contrast to term rewriting systems (Baader and Nipkow 1998) the definition of a query as two sets of triples and predicates simplifies the pattern matching, i.e. the identification of the right subpart of the query for the pattern match. A subset of both sets has to be determined which must be syntactically equal to the instantiated pattern. Please note that due to set semantics, the triples and predicates in the pattern may be distributed over the query. Now we will define how the new rewritten query is constructed with the help of the rewriting rule and pattern matching.

Definition 5 (Query Rewriting) If a rewriting rule $R = \langle PA, RE, CN \rangle$

- is applicable to a query $Q = \langle M_Q, O_Q, P_Q \rangle$ with subsets $M'_Q \subseteq M_Q, O'_Q \subseteq O_Q, P'_Q \subseteq P_Q$ and substitution θ
- $\theta(CN)$ is satisfied,

then the rewritten query $Q^R = \langle M_Q^R, O_Q^R, P_Q^R \rangle$ can be constructed with $M_Q^R = (M_Q \setminus M'_Q) \cup \theta(M_{RE}), O_Q^R = (O_Q \setminus O'_Q) \cup \theta(O_{RE})$ and $P_Q^R = (P_Q \setminus P'_Q) \cup \theta(P_{RE})$ with $RE = \langle M_{RE}, O_{RE}, P_{RE} \rangle$.

Informally speaking, if the pattern matches a query and the conditions are satisfied then the matched pattern is substituted by the replacement. Applying the above rewriting rule to the basic query we get the following refined query:

$$\begin{aligned}
 M_Q &= (\{(Resource, subject, Subject), \\
 &\quad (Resource, title, Title) \\
 &\quad (Resource, description, Description), \\
 &\quad (Resource, language, Language)\},) \\
 O_Q &= \emptyset \\
 P_Q &= \{like(Subject, "inferenceengines"), \\
 &\quad (Language = "de")\}
 \end{aligned} \tag{1}$$

According to the rewriting rule, the triple (*Resource, title, Subject*) which is matched by P_A of the rule is replaced by (*Resource, title, Subject*), (*Resource, language, Language*) and P_Q is extended by (*Language = "de"*). Note that the language preference of `user42` is German because of the satisfied constraint *languagePreference*(*user42, "de"*) and "de" means German. Similarly, the query can be rewritten further to add learner performance constraints. Note also that a similar process is applied in the query relaxation process. Rewriting rules will then contain, for example, a pattern on subject replaced with a pattern on title.

2.3 Using rewritings

In general, rewriting is a very powerful approach to manipulate over-constrained queries. By replacing parts of a query, we can realize four types of actions and subsequently arbitrary combinations of these actions:

- *Making Patterns optional*—this provides a query which considers a situation that some of the resources do not have all the metadata annotations expected. For example, some resources might not have a subject, hence, the corresponding part of the query has to be made optional. A query then gives also results where the particular predicate relaxed to an optional predicate does not occur;
- *Replacing Value*—this provides a query where a particular predicate value is replaced with another value or a variable. This can be useful to find resources that do not directly match the user interest, but are concerned with a broader or related topic. Taxonomies may be used to provide siblings, more general terms, and so on;
- *Replacing Patterns/Predicate*—this provides a query where a particular triple resp. predicate in restrictions is replaced by another triple resp. predicate. A domain knowledge is employed for this purpose. In the case of our example, if a subject query is not satisfied, it may be replaced by a title query with similarity measures;
- *Deleting Patterns/Predicate*—this provides a query where a particular predicate is deleted from a query completely.

As such, these operations are independent of the application domain and the user preferences. The connection to the domain and the user can be made using a set of predicates in the condition of the rewriting rules that link the manipulation to specific aspects of the domain and the user model. In the following, we discuss two predicates for including information about the domain and the user into the rewriting process.

- | | |
|----------------------------------|---|
| domain-preferred-over(X, Y) | This predicate indicates that due to the special nature of the domain a certain relation or value is a better choice for retrieving results than another one. We can use this to relax queries by replacing a highly preferred value by a less preferred one that is more likely to deliver a result even if this result may be less exact. |
| user-preferred-over(X, Y, U) | This predicate is defined in the context of a specific user U and indicates that the user considers a certain relation more important or prefers a certain value |

over another one. We can use this predicate to relax queries by replacing highly preferred values by less preferred ones or for deciding which of the predicates in a query can be relaxed more easily, because the user considers it less important.

In order to make the relaxation smooth, we consider these predicates to be non-transitive.⁵ The concrete implementation depends on the representation of the domain and the user profile. In the following, we discuss the implementation and use of an abstract user and domain model based on semantic web technologies and explain how rewriting conditions can be checked based on these models using an RDF query.

3 Environment, preference and domain model

In order to include knowledge about the domain of interest and the preferences of the user into the query relaxation process, we have designed a general scheme for representing relevant knowledge independent of a concrete application. This general scheme exploits the meta-modeling capabilities of RDF to define aspects of the world we can take into account in the rewriting process (compare Fig. 2).

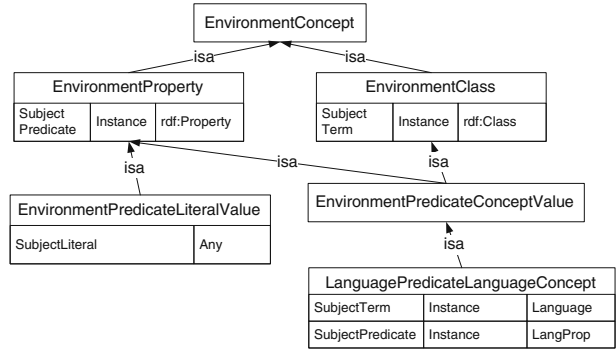
The schema follows an idea, that each environment can be generated according to an *application domain schema* used by the application. Rather than directly representing domain knowledge or user preferences it provides metaclasses that can be instantiated by existing representation schemes for information resources such as Learning Object Metadata (LOM) (Nilsson 2001) as well as metadata schemas like the Dublin Core standard (The Dublin Core Metadata Initiative 2008), and taxonomies and ontologies used for predicate values in the information resource schemas such as the ACM computing classification system (CCS) (Association of Computing Machinery 2002).

An environment concept can be, for example, linked to a field on a user interface form where the user can type a search term or it can be filled in with a class from a taxonomy. Such a generic environment schema provides us with the flexibility to describe any user environment which is based on schemas. For example, an environment concept can model a field on an entry form which is used to enter a subject term which a user is searching for in the metadata. Such a field will be an instance of the `EnvironmentProperty` class pointing to a “dc:subject” predicate of the Dublin Core schema. An example of combined class and predicate instances would be a predicate “dc:subject” with a class from a taxonomy like ACM CCS as its value.

Another advantage of such a generic environment schema is that we can refer to environment concepts from user preferences. Figure 3 depicts a schema for environment user preferences. Each user can express his level of preference for any environment concept. This is reflected by the `EnvironmentItem`

⁵Transitivity is ensured by the rewriting procedure itself; it does not need to be considered for the predicates themselves.

Fig. 2 A schema for generic environment



property of the `EnvironmentUserPreference` class. This is a generic definition of an environment preference through an `EnvironmentConcept` class as a domain for `EnvironmentItem` attributes. Classes for environment preferences are further specialized according to which environment concept class is used to describe them. For example, an environment class for representing language predicates `LanguagePredicateLanguageConcept` from Fig. 3 is used as a domain for the `EnvironmentItem` property of `LanguagePreference`. Note, that this class inherits and overloads properties from `EnvironmentProperty` as well as `EnvironmentClass`.

The level of a user preference can be expressed as a value from a metric. This is modeled by the `MesuredPreference` as a subclass of a user preference. The values from preference measures can be used to order them, i.e., to deduce the ordering relations between preference instances which is modeled by the `hasImportanceOver` relation. Another alternative is to deduce preference relations from usage logs as we describe in the next section.

Besides the user preferences, we also consider a schema for a user’s background. This is represented by the learning performance and the skills gained from learning. We use our schema for such a learner’s learning performance (Dolog and Schäfer 2005) where the learning performance is described by a relation to learning competences, portfolios created and certificates gained during/from learning activities which have been connected to the learning performance.

To show a concrete instance of the environment preferences of a user, let us now consider a situation where a user John whose id is user42, (cf. Fig. 1) prefers the

Fig. 3 A schema for environment user preferences

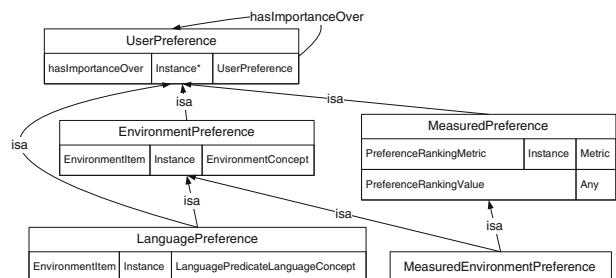
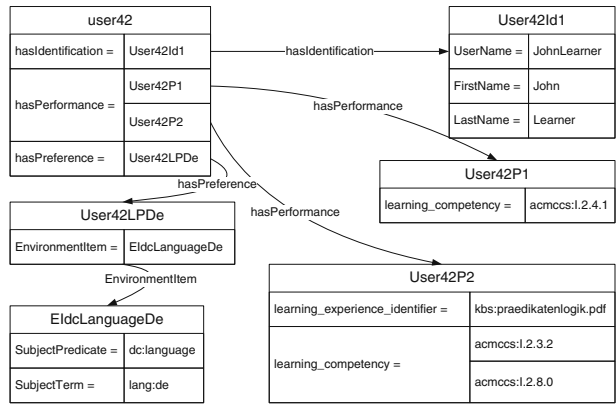


Fig. 4 An excerpt of instance examples for environment user preferences

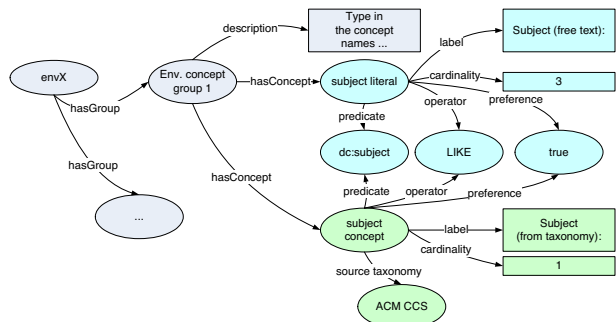


German language. In addition, he has attended two lectures, one on predicate logic and one on modal logic. An instance reflecting this situation described according to the environment user preference schemas is depicted in Fig. 4. His profile points to two performance objects: User1P1, and User1P2. The User1P1 is a performance record from a modal logic lecture where the user learned about the modal logic concepts (I.2.4.1 of ACM CCS), where User1P2 is a record from the predicate logic lecture where the user learned about inference engines (I.2.3.2 of ACM CCS) and backtracking (I.2.8.0 of ACM CCS). The user42 profile also points to one preference object: User1LPDe. This is a language preference referring to the German language (lang:de).

3.1 Preference based user interface and query generation

The user and environment preferences are utilized for the user interface and an initial query generation as well. Figure 5 shows an example of a group of related fields described in RDF and generated from the user preferences together with additional descriptions of those fields. There is a group of related fields for a user input on a learning goal (topic) field. There are two options a user always has: typing a free text and selecting from a taxonomy. This is reflected by two fields within the environment concept group. The free text can be typed in three distinct fields

Fig. 5 An excerpt of a user interface environment for query composition



(cardinality). Query related descriptions are predicate, operator, and preference. ‘Query predicate’ gives a restriction field to which a user input is being applied. ‘Operator’ (such as LIKE) gives a comparison operator to be applied for matching the user input with the predicate. ‘Preference attribute’ is given to specify whether a user can express his interest/preference utility for further query processing and query relaxation. Similarly to the free text subject field, a selection from the taxonomy field points to query-related attributes and a cardinality (in this case just ‘1’). Furthermore, it points to a default taxonomy to be used for concept generation.

A number of heuristics are used to produce such descriptions from the environment and user preferences. For example, the LIKE operator is usually generated in interface descriptions for text fields. A range operator is generated for time and date fields. For free text attributes, a triple of user interface fields is generated and one indented list is generated for taxonomical attributes. Taxonomies are suggested from pruning of resource data. If there is more than one taxonomy used in the resource data, a list with the used taxonomies is generated from which a user can select his preferred one. A user can store his own environment model with own preference values. There is always a default environment which is usually used for novice users.

This model is utilized in an initial query construction process as well. The descriptions like predicate, operator, and taxonomy are utilized when constructing the query restrictions. The preference values for attributes serve as an input to the relaxation process. Furthermore, the user preference model also contains information on projection predicates. Before the query is submitted to the relaxation service, other preferences stored in the user profile are considered as well for the initial query construction together with preference values (if available).

3.2 User preferences and user modeling

There are two major areas of related work for preference models: preferences in databases and CP-Nets. Pioneering work on preferences in databases and their formal models have been provided in Chomicky (2003), Kießling (2002). The preferences there are defined as partial orders over domains of attributes from a database schema. Preferences can be composed by different set operators such as union and intersection and so on. The composition creates another preference relation depending on the operator used to compose the preferences.

Another area where preferences have been considered is artificial intelligence. Conditional *ceteris paribus* represented as a network (Boutilier et al. 1999) or its extensions (Brafman et al. 2005) are used to describe preference dependencies over predicates. They describe how certain predicates with their values depend on the other predicates leading to different outcomes based on the preference predicates.

In our rewriting approach we combine benefits of both approaches. We allow for explicit representation of preference relation over RDF property domains. As in RDF schema properties are also defined as resources, we can consider similar relations over properties/predicates as well, thus providing a unified model describing preference relations over both predicates as well as domains. In the next section we show how both, preference relations over predicates as well as domains can be learned from RDF queries in existing systems.

In addition, we integrate such preference models with user profiling. User profiles contain preferences as its one dimension. Besides that, different aspects of

user context and history are useful in our context such as learning performance and experience. These aspects differ from preferences as they actually represent knowledge about past user actions relevant for the environment and a user task. Formally, they can form relations or nets similar to preferences but with different semantics. We also argue that preferences and user profiles strongly depend on the environment. Therefore, we provide a generic environment model where preferences can be embedded and related to schemas and ontologies of a problem domain provided through the environment.

Preferences in databases have been integrated with traditional query languages in preference queries. The preference queries were studied in the context of relational algebra (Chomicky 2003; Kießling 2002; Lacroix and Lavency 1987) or datalog (Kießling and Güntzer 1994; Köstler et al. 1995). The preference query languages provide several operators which are understood by engines implementing them. We adopt a different approach. As the combination of preferences and different aspects of user profiles depend on the context and environment, we rather argue for a more flexible approach, i.e., that rewriting rules, preferences, and user relevant dimensions are configured when an environment is deployed. Furthermore, the rules as well as the user profiles are continuously updated by learning preferences over selected domain and environment models. Our approach implements the preferences and user profiles by query rewriting at a level above a query language engine. This enables flexibility in the query language, the preference relations and the user profile as well as within the query engine.

Our approach relates to reasoning employed in CP-nets. However, while the main target of CP-net reasoning is to answer whether an outcome is preferred over another one based on the reasoning on *ceteris paribus*, our approach tries to improve the outcome by modifying a query.

4 Preference elicitation

An essential requirement for our approach to be useful in practice is the ability to deduce preferences as a basis for rewriting rules. As described above, we distinguish between user and domain preferences. The elicitation of user preferences is an active field of research as many personalization techniques rely on correct user models. We have positioned our work with respect to the major techniques proposed in the literature. Finding and representing domain preferences is a less well investigated problem and deserves more attention. In the following, we discuss the elicitation of domain preferences in the context of a concrete target domain. On top of that, there has been a proposal for query relaxation based on the RDF Schema data model, that can also be used as a generic preference model for applications where neither domain nor user preferences are available.

4.1 Domain preferences

In a recent article Hurtado and others have proposed a general query relaxation approach for RDF data based on the semantics of RDF schema (Hurtado et al. 2008). In particular, they propose a number of relaxation patterns for RDF query languages. These relaxation patterns are defined in terms of triple patterns and

can therefore directly be implemented in our approach. Besides general relaxation patterns for RDF that are similar to the general rewriting possibilities mentioned previously (dropping triples) and relaxations adopted from earlier research on databases (breaking join dependencies) the authors also propose ontology-based relaxation patterns that take the schema and therefore the nature of the domain into account. In Hurtado et al. (2006) the following schema-based relaxations are mentioned:

- Type relaxation: replacing a triple pattern (a, rdf:type, b) with (a, rdf:type, c), where (b, rdfs:subClassOf, c) follows from the model. For example, the triple pattern (?X, type, ConferenceArticle) can be relaxed to (?X, rdf:type, Article) and then to (?X, rdf:type, Publication).
- Predicate relaxation: replacing a triple pattern (a, p, b) with (a, q, c), where (p, rdfs:subPropertyOf, q) follows from the model. For example, the triple pattern (?X, proceedingsEditorOf, ?Y) can be relaxed to (?X, editorOf, ?Y) and then to (?X, contributorOf, ?Y).
- Predicate to domain relaxation: replacing a triple pattern (a, p, b) with (a, rdf:type, c), where (p, rdfs:domain, c) follows from the model. There are no domain declarations in Fig. 1.
- Predicate to range relaxation: replacing a triple pattern (a, p, b) with (b, rdf:type, c), where (p, rdfs:range, c) follows from the model. For example, the triple pattern (?X, editorOf, ?Y) can be relaxed to (?Y, rdf:type, Publication).

We consider these relaxations to be typical examples of domain preferences that can be used in our setting. The implementation of these patterns as rewriting rules is straightforward.

It turns out that following schema-based relaxation rules is often not enough in practical applications as not all relevant preferences have a direct relation to the schema. Useful relaxations often rather depend on the nature of user queries typically posed to the system. In order to get a better idea of the impact of the nature of user queries on domain preferences, we analyzed the problem of eliciting domain preferences for the REASE system (Diederich et al. 2007). REASE is a repository of learning resources for the domain of semantic web technologies. The system has been developed in the context of the Networks of Excellence KnowledgeWeb and REVERSE. It currently has more than 700 registered users. The system allows the user to pose keyword-based queries similar to search engines like Google. Internally, learning resources are represented by RDF metadata descriptions based on a complex schema. Details of the schema can be found in Brase (2005). This provides us with a certain degree of flexibility for mapping user queries onto the metadata description that we can exploit in the relaxation.

The REASE system already has a limited form of relaxation built into the search engine that has been designed based on extensive experiments for optimizing search results. In particular, each user query is evaluated against different metadata

Table 1 Variation of query terms in the REASE system (from Diederich et al. (2007))

Query term	Percentage	Hits
Problem solving method[s]	27%	2 / 2
psm[s]	22%	0 / 2
Problem solving methods psm	9%	2 / 2
Problem solving methods (advanced search)	7%	2 / 2
'Problem solving method[s]'	6%	1 / 2
Problem solving	5%	2 / 2

fields using different weights for computing the aggregated result. In particular, the following metadata fields are used as a target for user queries

1. Title (with weight 1)
2. Description (with weight 0.7)
3. MainContributorNames (with weight 0.3)
4. OtherContributorNames, EducationalObjectives, AdditionalInformation, Curriculum and Prerequisites: (with weight 0.1)

These empirically determined weights impose a preference relation over the different properties of a learning resource, that can be used in our system. Further, in the user study reported in Diederich et al. (2007) we analyzed the use of variations of search terms and the impact on the completeness of query results. Table 1 shows the result for the case of users searching for information about problem-solving methods. The analysis shows that users prefer the search terms 'problem solving methods' and 'psms' where the second search term will not return any result. We can use the information from the analysis and rewrite the query replacing the abbreviation 'psms' by the complete search term that is known to return results. Similar observations could be made for other search terms such as 'species' (searches for the different sublanguages of the web ontology language that are sometimes also referred to as 'layers'). This principle can be generalized by building a domain specific thesaurus as a basis for rewriting search terms.

In summary, since neither rewriting queries based on term lists and thesauri nor the schema-based rewriting of queries are new, the real benefit of our method is its ability to integrate different concrete approaches such as the ones mentioned into a common framework. In highly heterogeneous environments the combination of these different approaches is a real benefit that should not be underestimated.

5 Processing rewritings

In previous work (Dolog et al. 2006) we presented the theoretical foundations of rewriting RDF queries based on sets of triple patterns. In the following, we describe the concrete implementation of this rewriting approach in SWI PROLOG. We chose SWI PROLOG as a basis for the implementation because the declarative nature of PROLOG allows a straightforward implementation of rewriting systems. SWI PROLOG is especially suited because its semantic web library contains many useful functions for manipulating RDF data and RDF queries.

5.1 Rewriting queries

As a first step in the rewriting approach, the refined user query (cf. Fig. 1) is translated into a PROLOG representation. This is done by using the functionality provided by the PROLOG-based SeRQL query engine provided within SWI PROLOG. After this translation, the query mostly consists of a list of predicates of the form `rdf(Subject, Relation, Object)` that represent the FROM part of the query and a list of predicates of the form `serql_compare(Operation, Argument1, Argument2)` which represent the WHERE part of the query. These lists of predicates correspond to mandatory patterns⁶ and conditions in our rewriting approach mentioned in Section 2.1. Based on this representation, we can now also define rewriting rules in terms of special PROLOG predicates containing an identifier, a matching and a replacement pattern as well as conditions. The rewriting rules are applied to the PROLOG representation of the query. The resulting relaxed query is translated back into SeRQL using the corresponding functionality of SWI PROLOG.

The following rule is an example which makes use of the predicates introduced in Section 2.3. We can write generic rewriting rules that are guided by information from the environment and user model introduced in Section 3. Below we have an example of such a generic rule. This rule looks for relations in query patterns that are mentioned in the user model and replaces the corresponding relation name by a less preferred relation based on information about domain preferences with respect to relations.

```
'Generic-Domain-Preference' @@
    pattern(where([rdf(Subject, Relation, Object)]),
            from([]))
==> replace(where([rdf(Subject, Relation, Object),
                    rdf(Subject, LessPreferredRelation,
                        Object)]),
            from([]))
    && (domain_preferred_over(Relation,
                              LessPreferredRelation)).
```

This rule is used to solve the problem in example 1, where the preferred relation is 'subject' while the less preferred one is 'title'. The explicit representation of both, a user model and a meta-model of the domain presented in Section 3 provides us with a basis for computing the special predicates `user-preferred-over` and `domain-preferred-over` and using them to guide the query rewriting process. In particular, we can define these predicates in PROLOG using elements from the SWI RDF library to directly refer to the RDF-based representation of the user and environment model.

As the relations 'subject' and the relation 'title' are in the `user_preferred_over` relation with respect to our example user, the original query will be rewritten and the pattern `{Resource} subject {Subject}` will be replaced by `{Resource} title {Subject}`. A second rewriting rule can be used to modify the value of

⁶Currently optional patterns are not supported by the implementation, but the same machinery can be used to also include them into the rewriting.

‘Subject’ in such a way that the query verifies if the title contains the value of the ‘Subject’ variable as a substring. The combination of these two generic rewriting rules solve the problem in example 1.

5.2 Control strategy

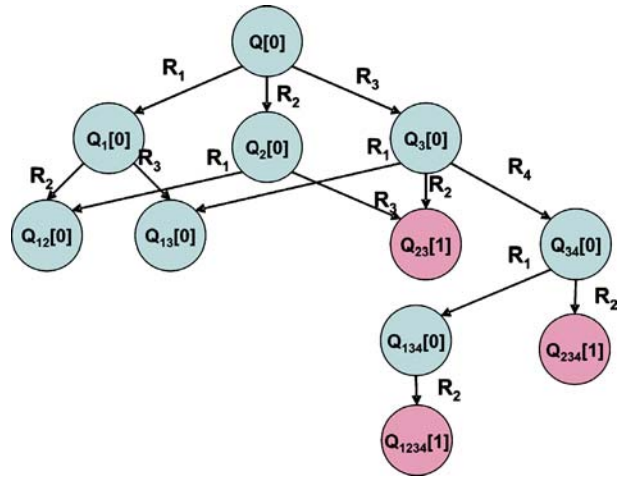
As mentioned above, the main problem of the rewriting approach to query relaxation is the definition of an appropriate control structure to determine in which order the individual rewriting rules are applied to queries. Different strategies can be applied to deal with the situation where multiple rewritings of a given query are possible. Examples are:

- User Interaction (Motro 1990): possible rewritings are presented to the user who decides in which direction to proceed
- Heuristic Search (Stojanovic 2003; Stuckenschmidt 2004): The best rewriting is determined based on the similarity of the resulting query with the original one.
- Divide and Conquer (i.e., Skylining) (Kießling and Köstler 2002; Lacroix and Lavency 1987): The best results of all possible combinations of rewritings are returned.

In the current version of the system we have implemented a simple version of skylining. In particular, we interpret the problem of finding relaxed queries as a classical search problem. The search space is defined by the set of all possible queries. Each application of a rewriting rule R on a query Q is a possible action denoted as $Q \xrightarrow{R} Q'$. A query represents a goal state in the search space if it does have answers. In the current implementation we use breath-first search for exploring this search space. Unlike classical search, however, the method does not stop when a goal state is reached. Instead, each goal state is explored and the results of the corresponding query are returned together with the relaxed query itself. As each goal state represents the best solution to the relaxation problem with respect to a certain combination of rewritings, the goal states form a skyline for the rewriting problem. The second difference to classical search is that we do not allow the same rule to be applied more than once with the same parameters in each branch of the search tree, because they only increase the complexity of query answering (Gutierrez et al. 2004).

Figure 6 illustrates the control strategy exploring the search space for the query given in Fig. 1. The nodes in the graph represent queries—the number in square brackets denotes the number of answers. The edges between nodes represent rewritings. There are four possible rewritings R_1 to R_4 , where R_1 relaxes the prerequisite requirement in the query, i.e. the first rewriting rule in Section 4.1. R_2 relaxes a string comparison in a SERQL query where two strings must no longer be equal but the first string can now be a substring of the second. R_3 and R_4 are two instantiations of the generic domain-preferred rewriting; R_3 replaces subject by title in the query whereas R_4 replaces title by description (cf. example 1). The search space is initialized by the original query Q . R_1 to R_3 are all applied resulting in the rewriting queries Q_1 to Q_3 . Because none of these rewritten queries return any results, the rewriting process proceeds with the application of R_1 to R_3 to the rewritten queries Q_1 to Q_3 . The resulting queries can be merged to Q_{12} , Q_{13} , and Q_{23} because it does not matter in which sequence R_1 and R_2 are applied. Additionally R_4 is now applicable to Q_3 resulting in Q_{34} because after replacing the subject by title (by R_3) the title can be

Fig. 6 Search space for example rewritings



replaced by the description. Furthermore, Q_{23} is a goal state because it returns one resource as an answer to the (relaxed) query. Figure 6 shows also how Q_{34} is further rewritten to the new goal states Q_{234} and Q_{1234} .⁷ Finally we get three rewritings Q_{23} , Q_{234} and Q_{1234} of the original query; each of them returns one learning resource.

6 Implementation

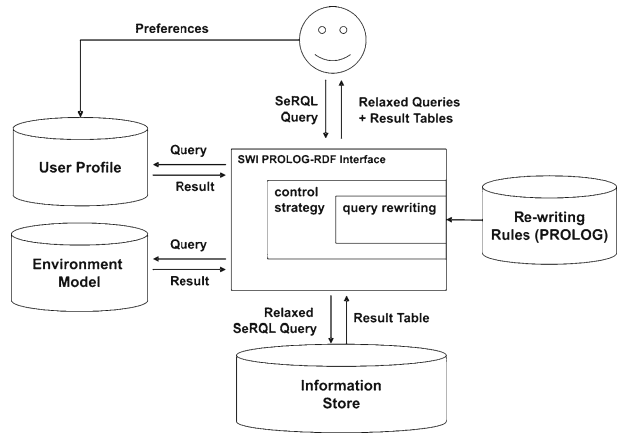
We have implemented the relaxation approach described above using state of the art semantic web technologies in the context of the European Research Networks KnowledgeWeb (Realizing the Semantic Web) and Prolearn (Network of Excellence in Professional Learning). The resulting system is an extension of the general e-learning infrastructure described in Dolog et al. (2004) which we verified on a repository of e-learning resources in the area of general computer science with the resources being annotated with RDF-based metadata using the schema proposed in Brase (2005). As the resources are provided by a large number of different authors, the metadata descriptions that form the basis of search contain many of the problems mentioned above which makes this data set an ideal test case for our approach. In the following, we provide an overview of the general system architecture and explain the prototypical search interface available on the web.

6.1 System architecture

Figure 7 shows the architecture of the implemented system. The system consists of a central component that contains the query rewriting functionality and the relaxation strategy. This component receives a user request in terms of a SeRQL query that has been generated by the user interface by automatically refining the user request based on known user preferences. This query is translated into a PROLOG representation

⁷The expansion of other states is omitted.

Fig. 7 System architecture



using the functionality provided by the SWI-PROLOG semantic web library. The rewriting rules which are specified in a PROLOG knowledge base are then applied to this representation of the user query. Most of the rewriting rules contain conditions that refer to user or environment preferences or both. The corresponding conditions are evaluated over the corresponding RDF models that contain the environment and user data, again making use of the PROLOG RDF interface provided by SWI PROLOG. The result of this rewriting is a series of more general queries that are translated back into SeRQL and are issued to a Sesame RDF repository that contains the actual data. The results of these queries are returned to the user interface together with the relaxed query in order to enable the user to understand the basis on which the results were achieved. In the following we discuss the different functional components of the system in more detail.

6.2 User interface

The prototypical search interface of the system combines user preference elicitation with a query formulation dialog. The original version of the personalized search described in Dolog et al. (2004, 2008) had just a query formulation for restrictions of the subject of resources. We have extended the user interface with generating environment based on the environment schema, a default environment for novice users, and a user preference elicitation. A user interface of such a personalization search environment is depicted in Fig. 8.

The default environment consists of items for specifying subject concepts, title, description, and language as query literals. Each of the attributes on the user interface has a preference elicitation slider. The slider is used to specify a value measuring an importance of a preference of a particular attribute to a user. Internally, these numerical values are translated into qualitative relations describing a total order of importance over the different aspects. With respect to Example 1 these relations specify in which order subject, title and description are considered as a source of information about the topic of a resource. A button for opening a dialog where a user can specify the value preferences and their order is provided where it is appropriate (e.g., the *Subject Values Preferences* button or the *Language Preferences* button). The

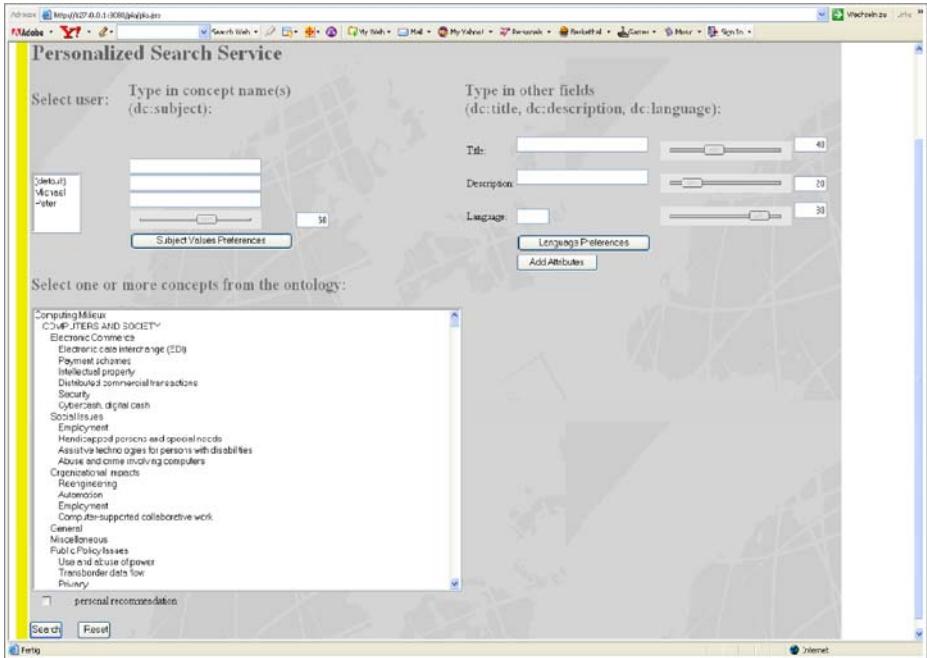


Fig. 8 Prototypical search interface

source of values for subject preferences is in our case the ACM CCS taxonomy, used also for selecting concepts on the user dialogs. We use the standard set of language identifiers as a source for values for the language preferences. The value preference dialog displays a tree, a graph or a set of concepts with value labels determining the importance of the preference. When a user points to a concept, a slide bar is drilled down to change the preference importance value. If a user needs to extend his restrictions, he can do that by selecting from other schema attributes which are offered when he presses the *Add Attributes* button. The attributes that a user filled in within the user interface are used to construct the restriction part of queries.

7 Conclusions

In this article we addressed the problem of querying RDF data containing irregularities due to multiple authorship and non-compliance to a standardized metadata schema. We illustrated this problem using an example from the e-learning domain, but we are convinced that the problem is a general one that is inherent in the idea of metadata on the semantic web. We have presented an approach for successively rewriting queries based on background knowledge. In particular, we have described how background knowledge about different alternative representations can be used to define generic relaxation rules that can be applied across different domains, provided that we have a suitable environment model. We have implemented and

tested the approach in the domain of e-learning using real world data about e-learning resources in computer science.

Open questions concerned with our approach are about suitable ways of acquiring the necessary information about user preferences as well as about the application environment. For the case of user preferences there is a large body of work in the area of user modeling including methods for automatically learning preferences based on user behavior. The acquisition of information about the environment model will probably be more of a challenge, because it is not clear whether alternative representations of the same information can be detected by observing the user. We have shown some strategies to acquire this information from usage logs. However, this problem needs to be studied further.

In summary, we can say that we need more experience with real data and real users in order to assess the effort connected with the acquisition of the knowledge necessary to successfully apply our approach in different domains. While for the domain of e-learning the benefits have been shown, this remains to be done in other domains.

References

- Association of Computing Machinery (2002). The ACM Computer Classification System. <http://www.acm.org/class/1998/>.
- Baader, F., & Nipkow, T. (1998). Term rewriting and all that. New York: Cambridge University Press.
- Boutilier, C., Brafman, R. I., Hoos, H. H., & Poole, D. (1999). Reasoning with conditional ceteris paribus preference statements. In *Proceedings of the Fifteenth Annual Conference on Uncertainty in Artificial Intelligence (UAI-99)* (pp. 71–80). Morgan Kaufmann: Stockholm.
- Brafman, R. I., Domshlak, C., Shimony, S. E., & Silver, Y. (2005). TCP-nets for preferences over sets. In *IJCAI-05 Multidisciplinary Workshop on Advances in Preference Handling at International Joint Conference on Artificial Intelligence*. Edinburgh, Scotland, (July). Available at: <http://wikix.ilog.fr/wiki/pub/Preference05/WebHome/P07.pdf>.
- Brase, J. (2005). Usage of metadata. Phd thesis, University of Hannover.
- Broeskstra, J., & Kampman, A. (2004). Serql: A second generation RDF query language. In *SWAD – Europe Workshop on Semantic Web Storage and Retrieval*. Amsterdam, The Netherlands, (November).
- Ceri, S. (1992). A declarative approach to active databases. In F. Golshani (Ed.), *ICDE* (pp. 452–456). Los Alamitos: IEEE Computer Society.
- Chomicky, J. (2003). Preference formulas in relational queries. *ACM Transactions on Database Systems*, 28(4):1–40 (December).
- Diederich, J., Džbor, M., & Maynard, D. (2007). REASE: The repository for learning units about the semantic web. *New Review of Hypermedia and Multimedia*, 13(2):211–237.
- Dolog, P., & Schäfer, M. (2005). A framework for browsing, manipulating and maintaining interoperable learner profiles. In L. Ardissono, P. Brna, & A. Mitrović (Eds.), *Proc. User Modeling 2005: 10th International Conference, UM 2005, LNAI*, (Vol. 2715). Edinburgh: Springer (July).
- Dolog, P., Henze, N., Nejdl, W., & Sintek, M. (2004). Personalization in distributed e-learning environments. In *Proc. of WWW2004 – The Thirteen International World Wide Web Conference*. New York: ACM Press (May).
- Dolog, P., Stuckenschmidt, H., & Wache, H. (2006). Robust query processing for personalized information access on the semantic web. In H. L. Larsen, G. Pasi, D. O. Arroyo, T. Andreassen, & H. Christiansen (Eds.), *Proceedings of 7th International Conference on Flexible Query Answering Systems (FQAS 2006)* (pp. 343–355), *Lecture Notes in Computer Science 4027*, 7–10 June 2006. Milan, Italy: Springer.
- Dolog, P., Simon, B., Klobucar, T., & Nejdl, W. (2008). Personalizing access to learning networks. *ACM Transactions on Internet Technologies. Special Issue on Distance Education*, 8(2) (May).

- Gaasterland, T., Godfrey, P., & Minker, J. (1992a). An overview of cooperative answering. *Journal of Intelligent Information Systems*, 1(2):123–157.
- Gaasterland, T., Godfrey, P., & Minker, J. (1992b). Relaxation as a platform for cooperative answering. *Journal of Intelligent Information Systems*, 1(3/4):293–321.
- Gutierrez, C., Hurtado, C., & Mendelzon, A. O. (2004). Foundations of semantic web databases. In *ACM Symposium on Principles of Database Systems (PODS)*. Paris, France, (June).
- Hayes, P. (2004). RDF Semantics. Recommendation, W3C.
- Hurtado, C., Poulouvasilis, A., & Wood, P. (2006). A relaxed approach to RDF querying. In *ISWC'2006 — 5th International Semantic Web Conference, Lecture Notes in Computer Science*. Athens: Springer-Verlag (November).
- Hurtado, C., Poulouvasilis, A., & Wood, P. (2008). Query relaxation in RDF. *Journal of Data Semantics*, 10, 31–61.
- Kießling, W. (2002). Foundations of preferences in database systems. In *Proceedings of 28th International Conference on Very Large Data Bases (VLDBO2)* (pp. 311–322).
- Kießling, W., & Güntzer, U. (1994). Database reasoning—a deductive framework for solving large and complex problems by means of subsumption. In *Proceedings of the 3rd Workshop on Information Systems and Artificial Intelligence, LNCS* (Vol. 777, pp. 118–138). New York: Springer.
- Kießling, W., & Köstler, G. (2002). Preference SQL - design, implementation, experiences. In *Proceedings of 28th International Conference on Very Large Data Bases (VLDBO2)* (pp. 990–1001).
- Köstler, G., Kießling, W., Thöne, H., & Güntzer, U. (1995). Fixpoint iteration with subsumption in deductive databases. *Journal of Intelligent Information Systems*, 4, 123–148.
- Lacroix, M., & Lavency, P. (1987). Preferences: Putting more knowledge into queries. In *Proceedings of the International Conference on Very Large Data Bases* (pp. 217–225). Amsterdam, The Netherlands.
- Motro, A. (1990). FLEXX: A tolerant and cooperative user interface to database. *IEEE Transactions on Knowledge and Data Engineering*, 2(2), 231–245.
- Nilsson, M. (2001) IMS Metadata RDF binding guide. <http://kmr.nada.kth.se/el/ims/metadata.html>, (May).
- Stojanovic, N. (2003). On analysing query ambiguity for query refinement: The librarian agent approach. In *Conceptual Modeling – ER 2003, volume 2813 of Lecture Notes in Computer Science* (pp. 490–505). Heidelberg: Springer-Verlag.
- Stuckenschmidt, H. (2004). Similarity-based query caching. In *6th International Conference on Flexible Query Answering System (FQAS), volume 3055 of Lecture Notes in Artificial Intelligence* (pp. 295–306). Lyon: Springer Verlag.
- Stuckenschmidt, H., van Harmelen, F., de Waard, A., Scerri, T., Bhogal, R., van Buel, J., et al. Exploring Large Document Repositories with RDF Technology: The DOPE Project. *IEEE Intelligent Systems*, 19(3), 34–40.
- The Dublin Core Metadata Initiative (2008). <http://dublincore.org/>.