



HAL
open science

Recommendation of XML Documents exploiting Quality Metadata and Views

Laure Berti-Équille

► **To cite this version:**

Laure Berti-Équille. Recommendation of XML Documents exploiting Quality Metadata and Views. 2nd Intl. Workshop on Data and Information Quality (DIQ 2005) in conjunction with the 17th Conference on Advanced Information Systems Engineering (CAiSE'05), Jun 2005, Porto, Portugal. pp.1-15. hal-01856346

HAL Id: hal-01856346

<https://inria.hal.science/hal-01856346>

Submitted on 10 Aug 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Recommendation of XML Documents exploiting Quality Metadata and Views

Laure Berti-Équille

IRISA, Campus Universitaire de Beaulieu
35042 Rennes cedex, France
Laure.Berti-Equille@irisa.fr

Abstract. In this paper, we propose to query XML documents with a quality-based recommendation of the results. The document quality is modeled as a set of (*criterion, value*) pairs collected in metadata sets, and are associated with the indexed XML documents. We implemented four basic operations to achieve quality recommendation: 1) annotation with metadata describing the documents quality, 2) indexing the documents, 3) matching queries and quality requirements, and 4) viewing the recommended parts of the documents. The quality requirements of each user are kept as individual quality profiles (called XPS files). Every XML document in the document database refers to a quality style sheets (called XQS files) which allow the specification of several matching strategies with rules that associate parts (sub-trees) of XML documents to user profile quality requirements. An algorithm is described for evaluation of the quality style sheets and user profiles in order to build an "adaptive quality view" of the retrieved XML document. The paper describes the general architecture of our quality-based recommender system for XML documents.

1 Introduction

Finding relevant, high-quality information in the World Wide Web or even in a collection of semi-structured documents is a difficult task. Information quality has no consensual definition and its evaluation requires: i) the measurement, ii) the weighted combination of both objective and subjective quality criteria and iii) the matching between the relative perception of information quality and the users' profile in terms of quality requirements.

The problem of high-quality information retrieval is becoming prevalent and not easy to solve today because of the growing, massive and quality-heterogeneous collections of documents now available on-line and also, because information quality in this context is very relative (depending on a topic, on a group of users, on a time period or on a focus of interest). As an introductory example, in intelligence gathering efforts for homeland security, information is collected from various sources with different degrees of trust and quality and then is corroborated (or not) by collaborating experts who need an automatic means to determine accurate and trustworthy information for decision making.

Content-based and collaborative recommender systems [RV97] work by automatically recognizing, tallying and redistributing recommendations of the web resources. The multi-confirmed recommendations appear to be significant resources for the relevant community and finally the number of distinct recommenders of a resource is a plausible measure of resource quality. But, many collaborative recommender systems particularly ratings-based systems are built on the assumption of role uniformity: they expect all users to do the same types of work in return for the same type of benefits. And the notion of user satisfaction and the evaluation task (rating) are very relative and should be considered in a flexible and adaptive way.

Our approach consists in taking into account the information quality evaluations and requirements in the context of collaborative annotation of XML documents for quality-driven information retrieval (IR) and information filtering (IF). We propose a modeling of document quality with various objective and subjective quality criteria. The objective criteria can be quantitative and calculated by statistical methods. The subjective criteria are defined and evaluated by a group of reviewers (or curators) collaborating on the qualitative annotation of the documents. The selection of the documents is both content-based and quality-based (*i.e.*, depending on the content and structure of the documents relevant to the query but also depending on both objective and subjective quality aspects required by the user). Our objective is to propose a multi-criteria adaptive recommendation of XML documents, and to refine the traditional selection of documents by exploiting embedded or linked metadata describing the resource quality. From these specifications, we developed the *XDARE* system (*XML-Documents Annotation and Recommendation Environment*) for quality-driven annotation and recommendation of XML documents.

1.1. Motivation

Among these various research propositions concerning on one hand, information quality, and on the other hand, recommender, blocking and adaptive hypermedia systems, we came to the point that there is no proposition in the literature nor system that combines resource quality for alternatively recommending/blocking and adapting digital resources on demand (in a way driven by users' profiles). Actually, in the current research works in IF or CF, the quality of the content is not considered as a key element for the users' decision in the recommendation process and we think this problem is not yet sufficiently addressed by existing approaches. Our motivation was to propose the three services (*i.e.*, blocking/recommending and adapting information) able to take into account in a flexible way the quality dimensions of the queried documents. Two principles guided our approach:

- in order to improve quality of a search result, it is necessary to evaluate the quality of the retrieved documents and to exploit it for the query processing,
- it's necessary that the definition of document quality remains flexible. The use of quality labels and metadata allow this flexibility for both specification and interrogation.

Compared to existing approaches for collaborative or content-based recommender systems, the innovative aspect of our approach is to include constraints and requisites on content quality that is complementary for better recommendation services.

1.2. Outline

The rest of the paper is organized as follows: section 2 presents the previous works on information quality and adaptive hypermedia and recommender systems. Section 3 describes our quality metadata model and presents the recommendation process for XML documents. Section 4 describes the architecture of our system. Lastly, Section 5 concludes the paper and presents our perspectives of research and development.

2 Related Works

2.1. Metadata and Data Quality

In the context of distributed information environments, metadata harvesting refers to the automatic collection of descriptive information from distributed resources. Recently, one particular way of accomplishing this collection of distributed metadata has been the subject of considerable attention in museums, archives and e-learning communities (e.g., the metadata collection proposed by the *Open Archives Initiative Metadata Harvesting Protocol* (OAI-MHP) [OAI02]). In the domains of geographical information systems [GJ98] and digital libraries (*Dublin Core, Bib-1, GILS, STARTS, Z39.50 ANS/NISO, etc.*), most of the exchange standards propose metadata specifications for information quality, which are either automatically extracted or measured by sampling from data sets. Many research works on information quality also proposed various definitions [WSF95], conceptual models [MP03] and methodologies [Wan02] [Red96] [AB+04] [SPP04] to improve or assess data quality in databases or in information systems [ICIQ96-04] [DQCIS03] [SN04]. The data quality dimensions most frequently mentioned in the literature are: *accuracy, completeness, actuality* and *consistency*. But many others dimensions [KSW02] [Nau02], metrics and measurement techniques [Win04][DJ03] have been proposed in the literature [LC02] [Red96] [FLR94] [Vas00] [MR00] [BP02] [Nau02] [NFL99]. Most of the techniques of quality measurement are centered on various methods of imputation such as inferring missing data from statistical patterns of available data, predicting accuracy estimations based on the given data, data editing (automating detection and handling of outliers in data), and error control. Concerning more specifically Web or semi-structured resources, the *DESIRE Project* [HBP00] also produced a detailed list of quality standards to be used for the selection of the Web resources with various categories of quality criteria: 1) criteria related on the policy of diffusion and the range of the resource, 2) criteria related to the content, 3) criteria related to the form, 4) criteria related to the management of the documentation quality.

In the context of the *TIPS European project* [TIPS99], several services have been developed related to the reuse of evaluations performed by humans on scientific publications. The first one, called *QCT (Quality Control Tools)* aims at collecting human detailed evaluations of documents in order to enrich the traditional topical indexing of documents with quality-related information (see Table 1 for the quality features used in the *QCT-TIPS* project for document quality). The second one, called *SF (Social Filtering)* integrates push functionalities as the alternate and complementary tool to traditional pull services such as information retrieval. Documents are pushed to users with respect to the evaluations they have made in the past, and compared to other users' evaluations.

2.2. Adaptive Hypermedia and Recommender Systems

A number of adaptive hypermedia systems have appeared as impersonalized systems, recommender systems with a common goal: to learn about the implicit preferences of individual users and to use this information to serve the entire community of these users better. The early recommender systems mainly used Information Filtering (IF) techniques and individual previous behavior to produce recommendations. To cope with the main drawback of IF techniques, Collaborative Filtering (CF) techniques have been proposed in order to recommend items based on the opinion (rating) of other users who have similar tastes. *GroupLens* [RI94] is a server-side recommendation engine for *Usenet* news. A user's profile is created by recording the user's explicit ratings of various articles. Automatic collaborative filtering is used to statistically compare one user's likes and dislikes with another user and to recommend articles from other similar users' profiles. Various recommender systems have been created to assist users for selecting potentially interesting information and for filtering out what users may not be interested in (such as *PHOAKS* [TH+97] for recommendation of Web resources mined from *Usenet* news messages, or *Ringo* [SM95], a music recommender system). But two major limits of the CF-based techniques are:

- the “*early-rater*” problem occurring for the first rating of documents without benefit of other previous recommendations,
- the “*sparsity rating*” problem occurring when the overlap between user's ratings (or number of co-rated items) is small or null and as a consequence that the recommendation results may be not accurate or cannot be produced.

The next level of recommender system is hybrid systems combining IF and CF techniques, such as *MovieLens* [GS+99], a movie recommender system using filterbots (IF agents) as rating robots which participate as members of the CF system. Both *Personal Web Watcher* [Mla96] are content-based systems that recommend web-page hyperlinks by comparing them with a history of previous pages visited by the user. *Personal Web Watcher* [Mla96] uses an offline period to generate a bag of words style profile for each of the pages visited during the previous browsing session. Hyperlinks on new pages can then be compared to this profile and graded accordingly. Most of the current hybrid systems still use co-rated items among users in finding correlated neighbors for an active user, and co-rated items between user and filterbot to find agreed filterbots.

On the opposite, PICS (*Platform for Internet Content Selection*) [RM96] is an infrastructure which associates labels to the contents of the documents available on the Internet in order to block the access for non-authorized users. Originally conceived to help the parents and the professors to control the navigation of their children/pupils on the Internet, it makes it possible to affect any criterion on the labels which the system interprets to authorize or block the access to the documents. Complementary to the recommendation, the advantage of this approach is that the structure of a document can be enriched by adding labels which define the conditions of viewability (blocking or full access). Another aspect of information personalization is to adapt web content to users' preferences and also to the variations of the client environment, so that web pages can be prepared suitable for the client. Adaptive hypermedia systems [BSS00] can learn about the implicit and explicit preferences of individual users and using this information to personalize information retrieval processes. In this context many adaptive hypermedia systems have been proposed, such as *OnlineAnywhere*¹, *SpyGlass*², *FastLane*³, *QuickWeb*⁴, *ProxyNet*⁵, *Digestor* [BS97], *TranSend* [FG+98], and *Mobiware* [AC+98]. However, most of them only make adaptation of the web content under special conditions due to the lack of structural information of HTML content, and many of them focus on image conversion.

3 Information Quality Metadata Modeling and Processing

3.1. Metadata for Document Quality

We proposed a simple XML document quality metadata schema (Figure 1) that: 1) allows a rigorous but flexible definition of each dimension of the quality of documents, 2) re-uses the existing standards of metadata proposed in the literature. We developed the corresponding system that 1) assists the user who can be a reviewer into the collaborative annotation process, and so, who should be able to define and evaluate himself the quality of the documents, 2) assists the user who's searching high quality information with providing him an adaptive recommendation (as quality view) of the retrieved documents depending both on his query and also on his quality requirements. The quality of a document is defined by combining quantitative measurements (computed by the system) and qualitative evaluation made by one (or several) reviewer(s); these metadata stored in XML files are embedded or linked to the content of the XML documents. As Figure 1 shows, the metadata type associated to a document (*metadatasetType*) can be a set of metadata (*metadataset*) or one metadata (*metadata*) which is composed of a criterion, a scoring value for the criterion given by a human reviewer (annotator) or computed by a program (*program*), a creation date and

¹ OnLineAnyWhere, FlashMap, <http://www.onlineanywhere.com>

² Spyglass, "White Paper of Prism 2.2", <http://www.spyglass.com/images/Prism22.pdf>

³ FastLane, <http://stage.acunet.net/spectrum/index.html>

⁴ QuickWeb, <http://www.intel.com/quickweb>

⁵ ProxiNet, ProxiWare, http://www.proxinet.com/products_n_serv/proxiware/

a comment (comment). A consensus can be calculated for a given criterion and a date if several notations have been proposed by several reviewers.

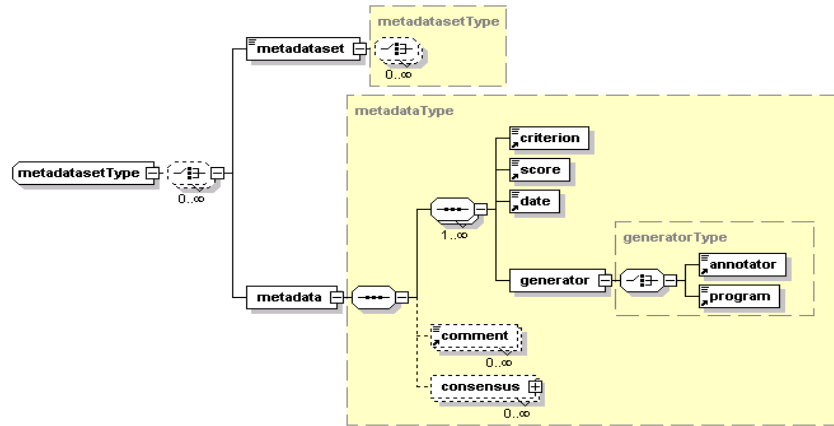


Figure 1. Metadata XML Schema Representation

Example 1. Figure 2 gives an example of quality metadata instances that can be associated to a document with both subjective criteria (originality, credibility) and objective criteria (citing popularity, reading popularity). In this example, the notations are values in $[0,10]$; the originality and the credibility of the document are evaluated by the reviewer A1 and the citing and reading popularity computed by programs similar to the one used in *CiteSeer* (<http://citeseer.nj.nec.com/>).

```

<metadataset qid="q1" type="ContentMD" scheme="" sortkey="" index="noindex" show="noshow" dldref="dld">
<metadata mid="m1">
  <critereion>Originality</critereion>
  <score>6</score>
  <date>12/02/2005</date>
  <generator><annotator>A1</annotator></generator>
</metadata>
<metadata mid="m2">
  <critereion>Credibility</critereion>
  <score>7</score>
  <date>15/02/2005</date>
  <generator><annotator>A1</annotator></generator>
</metadata>
<metadata mid="m3">
  <critereion>Citing Popularity</critereion>
  <score>7</score>
  <date>16/02/2005</date>
  <generator><program>Citation_Index_Program</program></generator>
</metadata>
<metadata mid="m4">
  <critereion>Reading Popularity</critereion>
  <score>6</score>
  <date>18/02/2005</date>
  <generator><program>Nb_Download_Program </program></generator>
</metadata>
</metadataset>

```

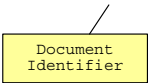


Figure 2. Example of quality metadata linked to the document identified by "dld"

3.2. Annotation and Recommendation Processing

Collaborative annotation and recommendation of the documents can be decomposed into four main operations (Figure 3): 1) harvesting and generating the quality metadata (annotating), 2) indexing the documents and their meta-descriptions (indexing), 3) matching the query (including the quality requirements of the users) with the representation of the indexed documents (matching), and 4) scoring and viewing the documents (viewing) with different recommendation strategies. Each step constitutes an elementary operation on the XML document collection. First, at the annotation step, each document of the collection is enriched by several quantitative and qualitative quality metadata (respectively, objective measures and subjective evaluations of the document quality) such as, in the example given in Figure 2. Metadata, document contents and structures are then used for the indexing process in order to represent each document as a multidimensional vector on these three axes: quality (criteria), content (terms) and structure (XML elements, attributes).

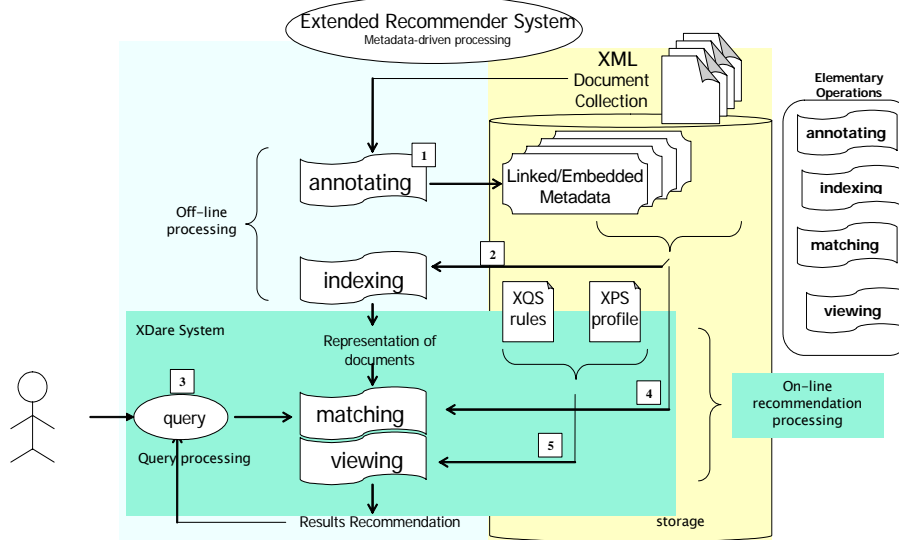


Figure 3. Operations for Retrieving Recommended XML Documents

For the query processing, the system finds the documents that best match the query terms and it applies to their respective metadata:

- the quality rules specified into the so-called XQS file (*XML Quality Style sheet*)
- the quality requirements specified into the so-called XPS file (*XML Profile Style sheet*) that correspond to the quality profile of the user (or of the group of users).

At this stage, the system has to use the recommendation strategy defined into the XQS file, to verify and to adapt (*i.e.*, soften) the quality constraints (given in the XPS profile) for building the quality-driven view of the query result.

Definition 1. XML Quality Style Sheet (XQS)

A quality style sheet (XQS) is an XML file used for the query processing. It references the profile to use as an attribute and is composed of the rules to apply for each document node. The DTD of a quality style sheet is given in Figure 4.


```

<!ELEMENT xqs (rule*)>
<!ATTLIST xqs
  DefaultProfileFile CDATA #FIXED "profiles.xps">
<!ELEMENT rule EMPTY>
<!ATTLIST rule
  Profile CDATA #REQUIRED
  DocumentNode CDATA #REQUIRED
  Priority CDATA #REQUIRED
  Access (default | recommending | blocking) #REQUIRED
  Strategy (default | exact | approximate | negotiated) #REQUIRED>

```

Figure 4. DTD for XQS sheet

Each rule is specifically applied to one document node with a specified user profile, priority (between 1-lowest and 5-highest), access and strategy attributes. The access defines the mode with or without recommendation or blocking access. The strategy defines how the quality constraints should be checked: exact matching, approximate matching or negotiation of the quality constraints. The exact strategy means that all the constraints must be verified on the value of each quality criterion (*i.e.*, metadata values) of the targeted document node. The approximate strategy allows the approximate and flexible matching between the constraints values and the effective quality criteria values using nearest neighbor algorithm on Euclidean distances between the multidimensional quality vectors. The negotiation strategy allows softening interactively the quality constraints in order to match the document nodes that answer the query with the best quality. The consistency of an XQS file is checked in order to avoid the definition of rules that give access (e.g. with recommendation) to the children nodes of a node whose access is forbidden (cf. *infra* Heuristics 2).

Example 2. Quality Style Sheet Example

In Figure 5, four rules are defined in the quality style sheet file concerning respectively three different users: the chief of the editorial board of a scientific journal, the secretary and the authors who submitted a paper for the special issue of the journal. The DTDs of the journal and the metadata set are given in Figure 7. These users will be allowed (or not) to access information according to the following rules:

- R1 – the secretary access all the submitted articles;
- R2 – the authors only access their own article;
- R3 – the authors are not allowed to access the papers and reviewers' comments of other authors;
- R4 – the editorial chief access the best submitted papers.

Definition 2. XML Profile Style Sheet (XPS)

A profile style sheet (XPS) is a file that references the user or the group of users of the profile and defines the constraints on the document quality dimensions. The set of quality constraints is defined as a quality contract.

```

<xqs DefaultProfileFile="profiles.xps">
  <!-- Rule 1 -->
  <rule Profile="users/secretary"
        DocumentNode="article[@id=$user]//node()"
        Priority="2" Access="default" Strategy="exact"/>
  <!-- Rule 2 -->
  <rule Profile="users/* [name()='Authors']"
        DocumentNode=" article[@dId=$user]//node()"
        Priority="5" Access="default" Strategy="exact"/>
  <!-- Rule 3 -->
  <rule Profile=" users/* [name()='Authors']"
        DocumentNode=" article[@dId=$user]/@qIdref//node() |
                      article[@dId!=$user]//node"
        Priority="5" Access="blocking" Strategy="default"/>
  <!-- Rule 4 -->
  <rule Profile="users/EditorialChief"
        DocumentNode="article//node() | article/@qIdref//node()"
        Priority="5" Access="recommending" Strategy="negotiated"/>
  <!-- Rule 5 -->
  <rule Profile="users/secretary"
        DocumentNode="article/@qIdref//node()"
        Priority="5" Access="blocking" Strategy="exact"/>
</xqs>

```

Figure 5. Example of XQS sheet

For a compact and simplified presentation, Figure 6 shows an extract of the XPS in the BNF-style grammar and an example corresponding to the user profile of the Editorial Chief considering the quality metadata given previously in Figure 2.

<pre> profile ::= PROFILE OF users { requisites } users ::= users member member member ::= user_name user_name ::= literal requisites ::= requisites requisite requisite ::= REQUIRE contractList contractList ::= contractList , contractElem contractElem contractElem ::= contractDefinition contractDefinition ::= CONTRACT { constraints } constraints ::= constraints constraint constraint constraint ::= dimName constraintOp dimValue dimName { aspects } constraintOp ::= == >= <= > < LIKE != dimName ::= literal dimValue ::= literal unit literal aspects ::= aspects aspect aspect ::= NUMBER constraintOp dimValue constraintOp dimValue freqRange constraintOp NUMBER % dimValue INIRangeLimit dimValue , dimValue rRangeLimit freqRange ::= [(lRangeLimit ::= [(rRangeLimit ::=]) </pre>	<pre> PROFILE OF EditorialChief { REQUIRE CONTRACT {Originality > 6 ; Accuracy > 6 ; Citing_Popularity > 7 per year ; Reading_Popularity in [6,9] per year } ; }; Quality Contract of the Editorial Chief </pre>
---	---

Figure 6. Extract of the XPS grammar and an example in the BNF-style

Example 3. Profile Style Sheet Example

The example given in Figure 6 presents the constraints required by the editorial chief on the document collection for what concerns the originality, the credibility, the citing and reading popularity of the articles. In particular, this user is interested in docu-

ments with a certain level of originality and accuracy (higher than 6 in the interval [0,10]), with citing popularity higher than 7 and reading popularity between 6 and 9. Figure 6 presents the user profile in a BNF-style grammar for simplification but, actually, the user profiles are stored in XML files using the RuleML⁶ DTD.

3.3. Recommending High-Quality XML Documents

The key elements of the recommendation processing are the quality-based recommendation rules (previously presented in XQS file) and the quality view that is generated according to two main heuristics we define in this section.

Definition 3. Quality-based Recommendation Rule

A quality-based recommendation rule is the following quadruplet:

Rule : < Profile, DocumentNode, Access, Strategy > *with*:

- *Profile*: the path expression related to the profiles of the XPS style sheet,
- *DocumentNode*: the Xpath expression evaluated inside the targeted XML document,
- *Access*: the value for recommending or for blocking the access to the information items (as document nodes),
- *Strategy*: the multi-criteria selection algorithm used to recommend the targeted node of the document. We use the following four methods for multicriteria selection defined and compared in [FC94]: *Weighted Linear Assignment (WLA)*, *Elimination aspect (EA)*, *Anderson method (AND)* [And90] and *Subramanian et Gershon (SG)*.

The specification of a quality-based recommendation rule is applied to node types (element, attribute, text...). Several integrity constraints may be defined in order to maintain consistency between the rules. Instead of recommendation of entire documents based on general quality requirements, we suggest the approach of filtering XML-document nodes based on quality criteria in the personal profile of the requestor. Our approach proposes document sub-trees recommendation that is used for complete document recommendation by aggregating quality scores of every document node.

Definition 4. Quality View

A quality view of a document is the result (as the fragment(s) of the XML document) that corresponds to the query and satisfies the quality constraints and rules defined in the user's profile. Two heuristics are used in order to build the quality view of each retrieved document in conformance with the quality constraints and the recommendation strategy chosen for the user who sent the query to the system. The quality view of the XML document is built node by node (i.e., XML element by element).

Heuristic 1. If the access to a node n of a document is allowed for recommendation for the user u , then u can see the recommended sub-tree of the XML document whose n is the root node if it satisfies the quality constraints defined as quality requisites into the user profile with the chosen strategy (exact or approximate or negotiated recommendation).

⁶ RuleML, <http://www.dfki.uni-kl.de/ruleml/>

Heuristic 2. If the access to a node n is forbidden for the user u , then u cannot access the sub-tree of the document whose n is the root node.

Quality views are then built according to:

- the quality-based recommendation rules given in the XQS file,
- the user's profiles given in the XPS file.

Example 4. Querying XML Documents and Quality Metadata

Following the Example 2, consider the edition of a scientific journal and the associated review comments of the submitted papers. The reviewers' comments are stored as metadata files. The content of the special issue of the journal has the following DTD (Figure 7):

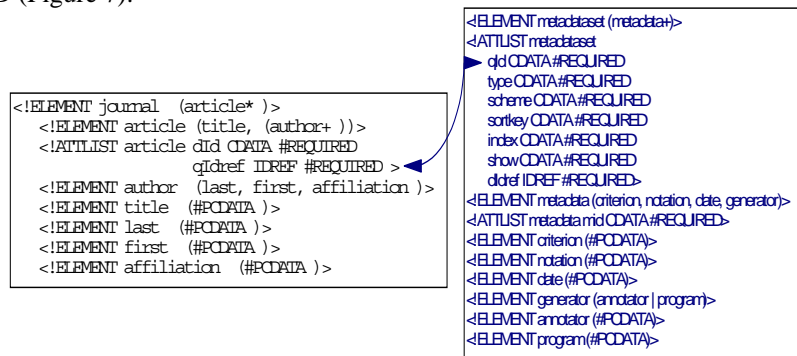


Figure 7. Example of DTDs for the journal and the reviewers' comments

Suppose the query given in Figure 8 sent by the three following users: the secretary, the authors of submitted papers and the editorial chief in order to get as a result XML element the name of authors, the title and the reviewers' comment of the submitted papers to the journal. The query language used is XQuery⁷ [FM01] [MM+01].

```

<results>
  {for $b in document /journal/article,
    $a in $b/author,
    $t in $b/title,
    $r in document($b/@qIdref)/metadataset
  return
  <result>
    { $a }
    { $t }
    { $r }
  }
</results>

```

Figure 8. Example of query

Then, the recommender system will return respectively the following results:

Result 1. Secretary's quality view. This query sent by the secretary will create as a result a flat list of all the author-title pairs of the submitted paper. Her quality requirements are given in the XQS file of Figure 5.

Result 2. Author's quality view. This query sent by an author will show only the title and the author name and the reviewers' comments of his own submitted article, but this author will not be allowed to see the papers and the comments of other authors.

⁷ XQuery, <http://www.w3.org/TR/xquery-semantics/>

Result 3. Editorial chief quality view. This query sent by the Editorial Chief will show the title, the author names and the reviewers' comment of the best articles.

4 System Architecture

We used standard tools for implementing the quality recommendation processor of our system called *XDARE (XML-Documents Annotation and Recommendation Environment)*. The description of the architecture is connected to the on-line processing steps showed in Figure 3 (numbered 3, 4 and 5). The prototype has been developed in Java and the following operations are implemented (Figure 9):

- ① **XML Document Analysis:** The document file is parsed (*Xerces Apache*) and syntactically analyzed (*SAX API*) and the corresponding internal representation is created and stored. The internal representation and manipulation has been developed with *DOM API*: the events produced by the *SAX API* during the first step of the document analysis are used to build the corresponding *DOM* tree. The *XDARE* operators use the instances of this internal representation by recopy and each query produces a new tree.
- ② **Query and Xpath Analysis:** The grammar is an extended XML Query syntax. *Jlex* and *Cup* are here used to produce the corresponding Java syntactical analyzer including the Xpath expression analysis. The query tree is explicitly instantiated for future optimization.

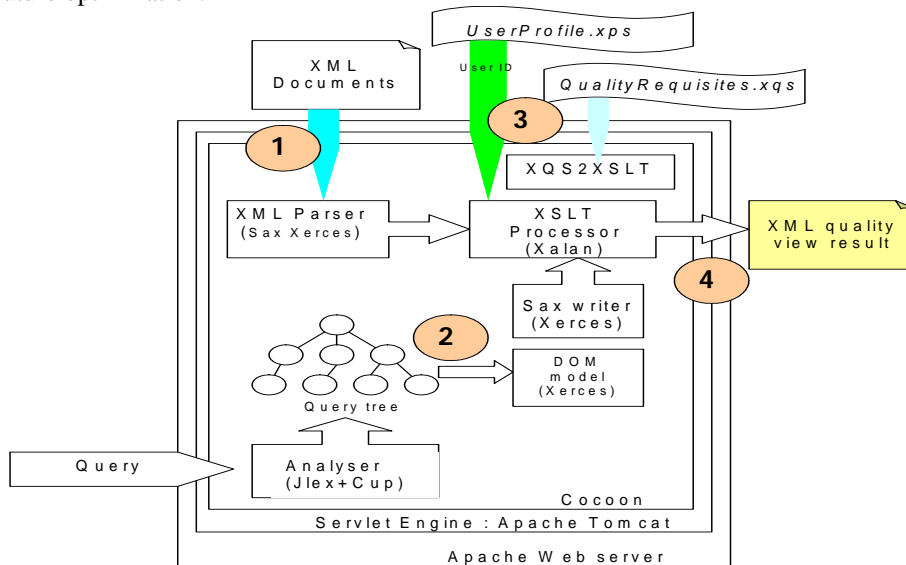


Figure 9. Quality Recommendation Processor of XDARE from the document processing perspective

- ③ **XML Document Quality Control:** In the Cocoon architecture, the *Xalan* processor applies a XSLT style sheet to the XML document. Our prototype transforms the XQS style sheet (including the quality requisites) into a set of XSLT templates. The application of XSLT style sheets enables the creation of quality views corresponding to user quality profile (defined into the XPS file).
- ④ **Quality View Generation:** When the user wants to browse a XML document, he actually obtains a view of this document in conformance with his quality requirements.

5 Conclusion and Perspectives

In this paper, we present a general architecture for quality-based recommendation of XML documents. The document quality is modeled as a set of (*criterion, value*) pairs collected in metadata sets, and are associated with XML documents.

We implemented four basic operations to achieve quality recommendation: 1) annotation with metadata describing the documents quality, 2) indexing the documents, 3) matching queries and quality requirements, and 4) viewing the recommended parts of the documents. The quality requirements of each user are kept as individual quality profiles (XPS files). Every XML document in the document database refers to a quality style sheets (XQS files) which allow for specification of several matching strategies and contain matching rules relating parts (sub-trees) of XML documents to user profiles. An algorithm is described for evaluation of the quality style sheets and user profiles in order to build an "adaptive quality view" of the retrieved XML document.

We sketched the architecture of the *XDARE* system, which implements the proposed data structures to support quality-based retrieval and adaptive quality views of XML documents. The specification and the exploitation of metadata describing the quality of documents can improve the system effectiveness for information searching and filtering including quality-driven recommendation.

The innovative aspect of our work is to propose the three services (blocking/recommending and adapting information) in a flexible way and to combine content-based and quality-based recommendation with considering the quality dimensions of queried XML documents.

References

- [And90] Anderson E., Choice models for evaluation and selection of software packages, *Journal of Management Information Systems*, vol. 6, p. 123-138, 1990
- [AB+04] Avenali A., Batini C., Bertolazzi P., Missier P., A Formulation of the Data Quality Optimization Problem, *Proc. of the Intl. CAiSE Workhop on Data and Information Quality (DIQ)*, Riga, 2004.
- [AC+98] Angin, O., Campbell, A.T., Kounavis, M. E., Liao, R.-F., The Mobeware Toolkit: Programmable support for adaptive mobile networking. *IEEE Personal Communication*, 5(4):32-43, 1998.

- [AT97] Asnicar F.A., Tasso C., ifWeb : a prototype of user model-based intelligent agent for document filtering and navigation in the Word Wide Web, Proc. of the 6th Intl. Conf. on User Modeling (UM'97), 1997.
- [Bal97] Balabanovic M. , An adaptive Web page recommendation service, Proc. of the 1st Intl. Conference on Autonomous Agents, 1997.
- [Bro80] Brodie M.L., Data quality in information systems, Information and Management, 3, 1980.
- [BP02] Ballou D.P., Pazer H., Modeling completeness versus consistency tradeoffs in information decision contexts, IEEE TKDE, 15(1):240-243, 2002.
- [BS97] BickMore T.W., Schilit B.N.: Digstor: Device-independent access to the World Wide Web, Proc. of the 6th International World Wide Web Conference, pages 655-663, 1997.
- [BSS00] Brusilovsky P., Stock O., Strapparava C. (Eds.): Adaptive Hypermedia and Adaptive Web-Based Systems, Proc. of the International AH'2000Conference, LNCS 1892, Trento, Italy, August 2000.
- [DJ03] Dasu T., Johnson T., Exploratory Data Mining and Data cleaning, Wiley, 2003.
- [DQCIS03] Proceedings of the Intl. Workshop on Data Quality in Cooperative Information Systems DQCIS'2003, Siena, Italy, 2003.
- [DR92] Delen G., D. Rijsenbrij D., The specification, engineering and measurement of information systems quality, Journal of Software Systems, 17, pages 205-217, 1992.
- [FG+98] Fox A., Gribble S.D., Chawathe Y., Brewer E.A.: Adapting to Network and Client Variation Using Infrastructural Proxies: Lessons and Perspectives. IEEE Personal Communication, 5(4):10-19, 1998.
- [ICIQ96-04] Proceedings of the Intl. Conf. On Information Quality, MIT, Boston, from 1996 to 2004.
- [FC94] Fritz C., Carter B., A classification and summary of software evaluation and selection methodologies, Technical Report, Mississippi State University, 1994
- [FLR94] Fox C., Levitin A., Redman T., The notion of data and its quality dimensions, Information Processing and Management, vol. 30, no. 1, 1994.
- [FM01] Fernandez M., Marsh J. XQuery 1.0 and XPath 2.0 Data Model. W3C Working Draft 2001. <http://www.w3.org/TR/query-datamodel/>, 2001.
- [GJ98] Goodchild M., Jeansoulin R. (Ed.), Data quality in geographic information : from error to uncertainty, Hermès, 1998.
- [GS+99] Good N., Schafer J., Konstan J., Borchers A., Sarwar B., Herlocker J., Riedl J., Combining Collaborative Filtering with Personal Agents for Better Recommendations. Proc. of the 1999 Conf. of the American Association of Artificial Intelligence (AAAI-99), pages 439-446, 1999.
- [HBP00] Hiom D., Belcher M., Place E., People Power and the Web: Building Quality Controlled Portals, TERENA Networking Conference, 2000. <http://www.desire.org/>
- [KSW02] Kahn B., Strong D., Wang R., Information Quality Benchmark: Product and Service Performance, Com. of the ACM, vol. 45, no.4, 2002.
- [KZ00] Koksalan M, Zions S. (Ed.), Multiple Criteria Decision Making in the New Millennium: Proc. of the 15th Intl. Conf. on Multiple Criteria Decision Making (MCDM), Ankara, Turkey, Lecture Notes in Economics and Mathematical Systems, 507, 2000.
- [LC02] Liu L., Chi L., Evolutionary Data Quality, Proc. of the Intl. Conf. Information Quality, 2002
- [Mla96] Mladenic, D., Personal WebWatcher: Implementation and Design, Technical Report IJS-DP-7472, Department of Intelligent Systems, J.Stefan Institute, Slovenia, 1996.
- [MGT+97]Malone T., Grant K., Turbak F., Brobst S., Cohen M., Intelligent information sharing systems, Com. of the ACM, 30(5):390-402, 1997.

- [MM+01] Malhotra A., Marsh J., Melton J., Robie J. (eds): XQuery 1.0 and XPath 2.0 Functions and Operators Version 1.0. W3C Working Draft 2001. Available at: www.w3.org/TR/xquery-operators/
- [MP03] Missier P., Batini C., A Multidimensional Model for Information Quality in CIS, Proc. Of the Intl. Conf. on Information Quality (ICIQ), 2003
- [Mou96] Moukas A., Amalthea : Information discovery and filtering using a multi-agent evolving ecosystem, Proc. of the Practical Application of Intelligent Agents and Multi-Agent Technologies (PAAM'96), 1996.
- [MR00] Mihaila G. A., Raschid L., and Vidal M.-E., Using quality of data metadata for source selection and ranking. In Proc. of the WebDB'00 Workshop, pages 93-98, 2000.
- [Nau02] Naumann F., Quality-Driven Query Answering for Integrated Information Systems, LNCS 2261, Springer 2002.
- [NLF99] Naumann F., Leser U., and Freytag J., Quality-driven integration of heterogeneous information systems. In Proc. of the 25th VLDB Conference, 1999.
- [OAI02] The Open Archives Initiative Protocol for Metadata Harvesting (Version 2.0), 2002. <http://www.openarchives.org/OAI/2.0/>
- [PMB97] Pazzani M., Muramatsu J., Billsus D., Syskill & Webert : identifying interesting Web sites, Proc. of the 13th National Conference on Artificial Intelligence, 1997.
- [Red96] Redman T., Data quality for the information age, Artech House Publishers, 1996.
- [RI94] Resnick P., Iacovou N., Suchak M., Bergstrom P., Riedl J. GroupLens: An Open Architecture for Collaborative Filtering of Netnews. Proc. of 1994 Conf. on Computer Supported Collaborative Work, pages 175-186, 1994.
- [RM96] Resnick P., Miller J., PICS: Internet access controls without censorship. Com. of the ACM, 39(10):87-93, 1996. <http://www.w3.org/PICS>
- [RV97] Resnick P., Varian H., Recommender systems. Com. of the ACM, 40(3):56-89, 1997.
- [SM95] Shardanand U., Maes P., Social Information Filtering: Algorithms for Automating "Word of Mouth". Proc. of ACM CHI'95 Conf. on Human Factors in Computing Systems, pages 210-217, 1995.
- [SN04] Scannapieco M., Naumann F. (Eds), 1st Intl. ACM SIGMOD Workshop on Information Quality in Information Systems (IQIS), 2004.
- [SPP04] Scannapieco M., Pernici B., Pierce E., IP-UML: A Methodology for Quality Improvement based on IP-MAP and UML. Advances in Management Information Systems – Information Quality Monograph (AMIS-IQ), Sharpe M.E., 2004.
- [TH+97] Terveen L., Hill W., Amento B., McDonald D., Creter J., PHOAKS: A System for Sharing Recommendations. Com. of the ACM, 40(3):59-62, 1997.
- [TIPS99] TIPS Documentation, Quality Control Tools User Requirements, V-Framework Programme IST-1999-10419, 1999. <http://tips.sissa.it>
- [Vas00] Vassiliadis P., Data Warehouse Modeling and Quality Issues, PhD thesis, Department of Electrical and Computer Engineering, University of Athens (Greece), 2000.
- [Wan02] Wang R., Journey to Data Quality, vol. 23 of Advances in Database Systems, Kluwer Academic Press, Boston, 2002.
- [Win04] Winkler W.E., Methods for Evaluating and Creating Data Quality, Information Systems, vol.29, no. 7, 2004.
- [WSF95] Wang R.Y., Storey V.C., Firth C.P., A framework for analysis of data quality research, IEEE TKDE, 7(4):623-638, 1995.