



Toward XML-Based Data Warehouse Architecture

Rami Rifaieh

R&D-TESSI INFORMATIQUE

72 bis Rue Bergson, 42000 Saint-Etienne, France
Email: rrifaieh@tessi2i.fr

Nabila Aïcha Benkat

Department of Informatics, LIRIS-INSa de Lyon
69621 Villeurbanne, FranceEmail: nabila.benharkat@if.insa-lyon.fr

ABSTRACT

XML enables data to migrate from relational databases and other sources into future applications. It integrates structured and unstructured data to present new application and knowledge management opportunities. By another way, data warehousing is an essential element of decision support, which has increasingly become a focus of the database industry. Meta-data contains data dictionary and repository; it describes warehousing process, data storage and information delivery. This paper defines how XML affects and changes the concept of the data warehouse. It describes an XML based tool for a structured, reusable and more efficient data warehouse (DW). Otherwise, it shows the ability to specify the warehousing tools with rules defined in the meta-data dictionary.

1 INTRODUCTION

Data warehouse system is a collection of data used for decision-making. The success of data warehouses implementation for business intelligence activities implies an increasing demand for new concepts and solutions [2]. It includes growth across platforms, tools, and applications. In all cases, the integration of data from heterogeneous sources is the essential stage of warehousing process. These sources could arise from existing legacy systems, which continue to generate data from mainframes computers through client server architecture. Moreover, disparate Web applications such as data sharing, real time data interchange, e-business, B2B activities and systems as ERP provide information that may be used to populate the company data warehouse. Then, data warehouse should be improved to support new applications such as real-time data warehouse and data web-house.

Furthermore, meta-data is an essential information that defines the what, where, how, and why about the used data. It can range from the conceptual overview of the real world to detailed physical specifications for the particular database management system. Meta-data can be used by automated tools (e.g. indexing robots) to improve the data interpretation / exploitation. Specifically, meta-data of business warehouse contains management rules [2]. They define which elements are supposed to be decisional and how they are calculated.

XML enables data to migrate from relational databases and other sources into new applications. It includes the ability to exchange data between application programs and browsers, between application programs and other application programs etc.

The purpose of this paper is to show how XML and its standards support data warehousing process. Therefore, XML could be used in all the construction process, especially in legacy data extraction, input transaction capture, cleansing procedures, direct storage of XML, and front-end information delivery. As we also deal with the evolution of meta-data, the XSLT component associated with XML documents helps us to support progressive meta-data for management rules. XML schema can provide another interesting dimension to XML text by defining datatypes. Moreover, an automatic generation of XSLT transformation could be a solution for personalization process based on meta-data description.

We will focus on showing where and how the extensible markup language can reduce process complexity. We try to show how using XSLT makes it easier to acquire and transform data. We discuss also the

automatic generation of XSLT from progressive meta-data, and show the usefulness of XML schema in the loading process.

The paper is organized as follows. In the next section, we will glance over DW and its construction process. The section 3 will detail the advantages of using XML for DWs. The section 4 will describe the design of XETL (XML-based Extraction, Transformation and Load tool). We will conclude our work in the section 5 with the future perspectives.

2 THE DATA WAREHOUSE SYSTEM

Data warehousing constitutes the background to enable business intelligence solution, which lets organizations access, analyze, and share information internally with employees and externally with customers, suppliers, and partners.

According to standard data warehouse architecture, the data warehouses systems include:

- ETL or warehousing tools: a set of tools which are responsible of preparing the data constituting the Warehouse database;
- Restitution tools: the diverse tools, which help the analysts to make their business decisions, and;
- Meta-data: it brings together the data about all the components inside the DW.

2.1 ETL Tools

ETL (Extraction, Transformation, and Load) represents the warehousing tools, or population tools. Warehousing tools have challenge to provide maintenance capability, availability, task management, and evolution support. Data integration and reuse possibilities are wide open but not yet very well realized. Although, some tools provide reused functions, these solutions still limited. Indeed, existing functions do not allow users to utilize an existing transformation plan and specify it with parameters to create a new data warehouse.

Maintenance process is not better; for example, if a user wants to change a SKU (stock keeping unit) number definition from five digits to seven. How many programs need to be changed to affect this enhancement? For the most of existing tools, in order to enable this operation a query has to be formulated into the data dictionary. Then, the user has to update all the concerning programs. In the population process, our suggestion is to combine the ETL tools with XML, since our concern is first of all the diversity of data sources, data targets.

2.2 Restitution Process

Often the information in data warehouses is published to a company's intranet web site. HTML is actually used to build these sites. Thus, restitution tools should provide the ability to generate HTML pages from warehouse database.

Nevertheless, delivering operation results from restitution model to mobile portals is going to be a new feature for restitution tools. Wireless devices (Pocket and mobile phone) will be able to capture information and use it in a real time with geographically distant warehouse. Thus, new standard output format should be used such as WML (Wireless Markup Language).

2.3 Meta-data Management

Traditionally, meta-data is data about the data. The meta-data allows the various function areas in the warehouse to communicate. It encompasses all corporate resources: database catalogs, data dictionary, and data models [9]. In the common data architecture, data catalogs are needed to understand data sources. In general, the meta-data should cover the acquisition, access, and distribution of warehouse data and should be the key to provide the business user with a complete map of the data warehouse.

3 XML AND DATA WAREHOUSE SYSTEM

For information systems, XML includes the ability to exchange data between application programs and browsers, between application programs and other application programs etc.

3.1 How Can XML Improve Data Warehousing

XML can contribute in warehousing process: data integration, cleansing procedures, data storage, and front-end information delivery.

3.1.1 Source Data Integration

The integration of data from different sources is the essential stage of warehousing process. These sources could arise from existing legacy systems, which continue to generate data from mainframe computers through client server architecture. Moreover, disparate Web applications such as data sharing, real time data interchange, e-business, B2B activities and other systems such as ERP provide information that should be used to populate the company data warehouse.

If the systems submit data in XML format, then partners sharing a common XML schema can transmit and receive data. Although older legacy applications may need to be retrofitted with XML writers, many more modern systems are equipped to write data in XML format. Existing relational databases already support query output directly in XML form.

The ability to read remote data via XML greatly simplifies data extraction, because custom parsers for the remote data source do not have to be written. This, in turn, promotes a more distributed and opportunistic approach to a spread-out "data web-house". Data web-house is a data warehouse where sources data come from XML web documents; it constitutes an essential search engine for a collection of web data.

3.1.2 Native XML Data Storage

Native XML database [7] can be used to store the core of XML-based data warehouse. They provide the possibility to integrate semi-structured data inside the data warehouse. This direct XML storage can fulfil the needs of storage support for data web-house. In this category, we can identify *Tamino* (<http://www.softwareag.com/tamino/>), *Natix* (<http://www.dataexmachina.de/natix.html>), and *Lore* (<http://www-db.stanford.edu/lore/>). Hybrid scenarios in which data in XML form could be joined inside the database directly to data in conventional relational tables is also interesting. Indeed, market DBMS offer the capability to bulk load data inside relational table from XML sources. This step is done with the help of XML schema. Furthermore, an XML data type in a relational system [4] would open up completely new classes of hierarchical applications.

3.1.3 Front-end Information Delivery

More and more sophisticated tools are used to deliver information outside the data warehouse. The widespread deployment of XML will be helpful in removing query and reporting tools from end users' terminals. An XML data transfer plus an associated XSLT formatting specification is enough to produce any desired user interface presentation on a remote browser. The result business information could be published with HTML pages in the company's intranet web site.

Moreover, different end-user clients can use different XSLT formatting specifications while drawing data from the same database server. All that the server needs to do is expose its query results through XML and from the client side just a Web browser is needed.

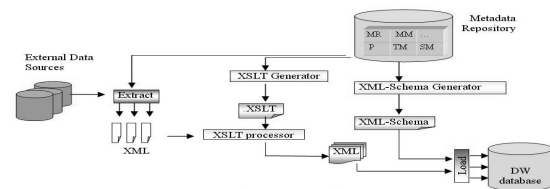


Fig 1. XETL Architecture

3.2 Advanced Features

3.2.1 Efficiency

XML is a neutral format, which can be used to enable the communication between different warehousing tools. Indeed, no owner format is needed to link and merge different warehousing tools, which is useful to divide the process and use a tool for extraction another one for transformation etc. The efficiency of the global process will be higher, especially when tools have good performances within different steps. We can chain a high extraction tool with XML generation capability and a high transformation tool with XML input, etc. Since XML is used in all the tools, no extra transformation between different owner tools formats is required.

3.2.2 Maintenance

The business rules of an organization are never the same, due to changes in the real world. So far, how do we maintain consistency when business rules change as a result of corporate reorganizations, regulatory changes, or other changes in business practices? How many places are impacted for each of these potential changes?

Due to the evolution of management rules, many programs need to be modified to enable new decision elements or new formulae for restitution. These rules define the decisional elements, equations, and parameters. For the most of existing tools, in order to realize this operation a query has to be formulated into the meta-data. Then, the user has to update all the concerning programs. To make maintenance process work better, the warehousing process has to be updated automatically. XML presents the ability to specify the transformations using XSLT. If the system is able generate automatically a new transformation, no extra update will be needed to enable the evolution.

3.2.3 Scalability

An important aspect to achieve the scalability is that changes in the implementation on either the application system side or the data side are not changing the transported data format between the two applications. This means that the client application or the server application does not need to be aware of that change.

One way of achieving that is to move the data in the standardized format between these components. Currently, the main format that is used in all the scenarios is XML. The idea is to take the data, in the data system or in the application system, move it into the standardized format of XML, both on the wire as well as using standardized grammar, which you can express in XML, to warp the data and transport it between the different systems.

3.2.4 Interoperability

Corporations are interacting on the web, blending suppliers, partners, and customers to form a virtual enterprise that function as the superset of the physical organizations. This e-business is performed through a myriad of real time information exchange technologies such as electronic document interchange (EDI), business-to-business exchange and e-business server applications.

Establish interoperability between these different sources is needed to integrate data inside DW. XML can play a major role in the integration of web sources inside the business DW, from corporate information portals, product catalogs, parts databases, and business-to-business document exchange etc.

3.3 How Meta-data Improve Data Transformation

For actual data warehouse architecture, few of tools communicate directly with meta-data repository. Meta-data are used to answer administrator queries. These queries offer information concerning the structure, models, and the warehousing process. Meanwhile, the set of meta-data is passive and query limited.

The existing tools do not deal with the evolution of management rules and reusability. Indeed, by one way, meta-data of business warehouse contains management rules. A passive use of meta-data implies two sorts of updates: meta-data repository and programs that handle data into warehouse. By another way, ERP systems [3] store data of different applications by a similar way on mainframe computers. In this case, data dictionary should provide more than just data attribute descriptions, range value, valid value, etc. It should supply the personalization process by applying the needed mapping to achieve the construction of different data warehouse. Thus, the same data element might be used by different entities of different applications to mean different things. Different data element names could also be used to represent the same things, potentially creating hundreds of instances of the same data all inconsistently named. In this case, data dictionary should provide more than just data attribute descriptions, range value, valid value, etc. It should supply the personalization process by applying the needed mapping to achieve the construction of different data warehouse. A generic plan is needed to establish mapping between source and target data. If meta-data contains all the generic plans for a target application, a personalization process can be used to extract data from source and apply the needed mapping.

The idea of making meta-data being *active* can improve warehouse systems. The solution can be performed with traditional transformation query with an automatic query generator or it can be a specific transformation language as used in XETL tool.

4PRESENTATION OF XETL TOOL

In this section, we present XETL, an XML based Extraction, Transformation and Load tool (Fig.1). The basic idea of our tool consists of using the XML format to create new generation of data warehouse, where XML is used as a pivot format to perform the warehousing process. A full case study of using XETL with real scale commercial data is described in [12]. This study shows the effectiveness of using XETL as a warehousing tool to populate a relational database.

4.1 XETL and Existing Tools

XETL takes advantage of XML ability to realize interoperability, scalability, maintenance, efficiency, and data integration opportunity. We consider that data transformation process is generated directly from meta-data. Thus, the designer will not be called to verify the consistency of the process, since this process uses parameters extracted from meta-data repository. At the same time, traditional request over meta-data is accessible for system's administrator. As we are dealing with XML format, we have tried to perform this functionality by generating XSLT transformation from meta-data.

4.2 XETL Architecture

XETL's architecture (Fig.1) integrates data from different and heterogeneous sources. Moreover, XETL is an active ETL; it interacts with meta-data to give parameters for extraction programs, to generate the transformation, and to specify the loading with target schema. Hence, it optimizes the flow of data, reduces the update of warehousing process by making automatic the creation of valid transformation, and generates schema from meta-data repository. The different components of XETL tool are described below:

4.2.1 The Meta-Data Components

The Meta-data repository is the essential element in the system. It includes a set of information of which:

- The mapping model (MM): this model is used to describe the mapping expressions. By mapping expressions, we mean the needed information to identify how a target field could be mapped from a set of source

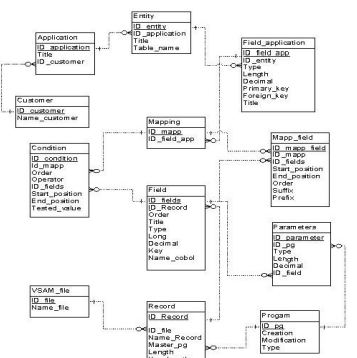


Fig 2. A mapping Meta-data model

fields. Fig.2 represents an example of MM used in [12].

- The source model (SM) contains the model of the source data. It covers relational, object oriented and semi-structured data modelling [11].
- The target model (TM) is alike to (SM); it describes the target data model. Moreover, it could cover the multidimensional model used by OLAP processing.
- Management rules (MR) is the set of rules defined by the administrator in order to fulfil business requirements.

4.2.2 The Extraction Process

The role of these programs consists in converting to a common XML format, the different sources. These sources could be traditional databases or inner digital documents that are produced by applications in enterprise or even web documents published by other partners. This first stage supplies the extraction program with the needed parameters from meta-data repository.

4.2.3 The XSLT Generator

The XSLT generator is a module, which can read useful parameters, rules, and mapping specification from meta-data repository to create a style sheet transformation. Concretely, a query is performed on the mapping model (MM). The result file is used by the Generator to produce the style sheet file. Then, an XSLT processor permits to execute the transformation on XML source documents. It generates a new collection of clean XML documents.

Indeed, the selection and filters get rid of superfluous data. During this process, all control and check are applied to data. Thus, the XSLT is more than cosmetic change for XML data; it tackles the content, structure and valid values. The architecture is not limited by the using of XSLT language; any other XML query language can be used to establish such process. For example, X-Query can be useful because it is more adapted for documents databases query and it is more optimized. On the other hand, the stability of XSLT (1.0) and the multi implementation is an advantage over the X-Query.

ModelTitle	TYPE	Table	Source	CONDITION	TITLE	POSITION
ID_customer	char(10)	Header_after_join_customer	element set	TOP3 = 'F' AND DOWNEERD='CAV' AND CODEINTER='CLIENT'	CLEMENT	
ID_HASC	char(5)	Header_after_join_customer	element set	TOP3 = 'F' AND DOWNEERD='CAV' AND CODEINTER='CLIENT'	CLEMENTE	
Day_HASC	YYYYMM DD	Header_after_join_customer	element set	TOP3 = 'F' AND DOWNEERD='CAV' AND CODEINTER='CLIENT'	DATE	
ID_LASC	num(10)	Line_after_join_customer	element set	TOP3 = 'L'	CLEMENT	1-10
ID_HASC	num(5)	Line_after_join_customer	element set	TOP3 = 'L'	CLEMENT	11-14
Type	char(10)	Line_after_join_customer	element set	TOP3 = 'L'	DOWNEERD	
Description	char(14)	Line_after_join_customer	element set	TOP3 = 'L'	DOWNEERD	1-14
Characteristic	char(30)	Line_after_join_customer	element set	TOP3 = 'L'	LIBELER	1-30
Line_status	char(2)	Line_after_join_customer	element set	TOP3 = 'L'	DOWNEERD	1-2
Weight	num(5,2)	Line_after_join_customer	element set	TOP3 = 'L'	DOWNEERD	3-7
Location	char(10)	Line_after_join_customer	element set	TOP3 = 'L'	DOWNEERD	

Fig 3. The result of query into mapping meta-data (MM)

A *XSLT-Generator* prototype written in C++ permits to create formatted XSLT transformations, this prototype is described with a case study in [12]. The *XSLT-Generator* communicates with meta-data to read needed mapping. Indeed, a query is performed on MM, which is a part of the meta-data, and the result Fig.3 is used to generate the transformations Fig.4.

The part of the meta-data used to produce the XML Schemas is the Target Model (TM). The generator reads the target model and creates an XML schema. In particular, to each entity of the target model is associated an XML schema. The most interesting point in this generation is that it supplies the loading process by providing schemas that include data types.

The fourth step consists on loading data to warehouse database with schema description. The result XML documents should be valid to XML schema. The loader deals with the validation of these documents. If an error occurs, the loading process is interrupted, and the error is searched for upstream. Therefore, the integrity of target DW database will be preserved.

Our work differs from the work done in Xyleme project [10]. Xyleme is XML base web-house, it stores XML pages in order to constitute an XML data repository. Xyleme deals with storage of huge quantities of XML data, query processing, and data acquisition strategies to update repositories. It provides control with services such as query subscription and semantic data integration to free users from having to deal with many specific DTD when expressing queries [1]. Our approach deals with different problematic. We are not studying the issue of building a data warehouse for the web (data web-house). We discussed the issue how XML can change the existing ETL tools and the advantage of such solution.

Copyright © 2003, Idea Group Inc. Copying or distributing in print or electronic forms without written permission of Idea Group Inc. is prohibited.

0 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

www.igi-global.com/proceeding-paper/toward-xml-based-data-warehouse/32072

Related Content

Explaining and Predicting Users' Continuance Usage Intention Toward E-Filing Utilizing Technology Continuance Theory

Santhanamery T. and T. Ramayah (2018). *Encyclopedia of Information Science and Technology, Fourth Edition* (pp. 831-846).

www.irma-international.org/chapter/explaining-and-predicting-users-continuance-usage-intention-toward-e-filing-utilizing-technology-continuance-theory/183796

An Introduction to Clustering Algorithms in Big Data

Rajit Nair and Amit Bhagat (2021). *Encyclopedia of Information Science and Technology, Fifth Edition* (pp. 559-576).

www.irma-international.org/chapter/an-introduction-to-clustering-algorithms-in-big-data/260214

Interpretable Image Recognition Models for Big Data With Prototypes and Uncertainty

Jingqi Wang (2023). *International Journal of Information Technologies and Systems Approach* (pp. 1-15).

www.irma-international.org/article/interpretable-image-recognition-models-for-big-data-with-prototypes-and-uncertainty/318122

Integrated Digital Health Systems Design: A Service-Oriented Soft Systems Methodology

Wullianallur Raghupathi and Amjad Umar (2009). *International Journal of Information Technologies and Systems Approach* (pp. 15-33).

www.irma-international.org/article/integrated-digital-health-systems-design/4024

Logistics Distribution Route Optimization With Time Windows Based on Multi-Agent Deep Reinforcement Learning

Fahong Yu, Meijia Chen, Xiaoyun Xia, Dongping Zhu, Qiang Peng and Kuibiao Deng (2024). *International Journal of Information Technologies and Systems Approach* (pp. 1-23).

www.irma-international.org/article/logistics-distribution-route-optimization-with-time-windows-based-on-multi-agent-deep-reinforcement-learning/342084