# MD2 – Getting Users Involved in the Development of Data Warehouse Applications

Gilmar M. Freitas
Prodabel
Av. Presidente Carlos Luz, 1275
31230-901 Belo Horizonte  MG
Brazil
gilmar@pbh.gov.br

Alberto H. F. Laender
Department of Computer Science
Federal University of Minas Gerais
31270-901 Belo Horizonte MG
Brazil
laender@dcc.ufmg.br

Maria Luiza Campos
Department of Computer Science
Federal University of Rio de Janeiro
21941-590 Rio de Janeiro RJ
Brazil
mluiza@nce.ufrj.br

## Abstract

User participation has become paramount for any successful data warehouse initiative. In this paper, we present MD2, a tool based on the dimensional data modeling approach, which facilitates the user participation in the development of a data warehouse application. MD2 assists users on identifying their analytical needs in order to help data warehouse designers to better specify business requirements and latter translate them into appropriate design elements. The tool is supported by a data repository that helps gather metadata used to build the application dimensional schema and to specify reports that can be generated from the data warehouse and visualized through a Web interface.

## 1   Introduction

Data warehousing has become a key technology for many organizations. From a general point of view, a data warehouse can be seen as an integrated data repository that is generated and then used across an organization for supporting the processes of strategic planning and decision making [4,7,8,10,18]. However, it is a general understanding that this is a complex technology that involves a number of difficulties for its effective implementation. Among the main difficulties pointed out by several authors [4,7,10,19] is the existence in the organizations of a large spectrum of users with distinct characteristics and information needs, which makes user participation paramount for any successful data warehouse initiative.

Dimensional modeling has been the predominant approach to designing data warehouse applications [2,3,10,11,15,18]. The main dimensional modeling primitives are *facts* and *dimensions*.

Facts represent measurements of the business or record events. Dimensions are attributes used to constraint, group or browse facts. On a relational implementation these primitives are represented by interrelated tables whose structure is usually described  by  a star schema [6,10,12,15]. Thus, the correct understanding and identification of facts, dimensions and their interrelationships are key to precisely design a data warehouse that conforms to the users' data requirements and analytical needs. To achieve this, it is very important an effective user participation when developing a data warehouse application.

In this paper, we present MD2, a tool based on the dimensional data modeling approach, which facilitates the user participation in several steps involved in the development of a data warehouse application. MD2 assists users in defining their analytical needs and modeling the application. The tool is supported by a data repository that helps gather metadata used to build the application dimensional schema and to specify reports that can be generated from the data warehouse and visualized through a Web interface. In addition to describing the main facilities provided by our tool, we also discuss, based on the Business Dimensional Lifecycle [10], a case study that illustrates how the tool is used to develop a data warehouse application.

The remainder of the paper is organized as follows. Section 2 discusses related work. Section 3 presents an overview of MD2. Section 4 discusses the interaction of MD2 with the Business Dimensional Lifecycle methodology. The case study is described in Section 5. Finally, Section 6 presents our conclusions.

## 2 Related Work

The concepts and techniques related to data warehouses have been widely discussed in the literature [4,7,10,11]. In what follows, we focus our attention on the dimensional model and on methodologies that adopt this approach (or some of its variations) to design, develop, and deploy data warehouse applications.

The dimensional model has been extensively addressed in the literature. A method for deriving dimensional schemas from traditional entity-relationship schemas is presented in [15]. The dimensional fact model, proposed by Golfarelli et al. [5], provides a graphical notation for representing data warehouses based on the dimensional model. The mapping of ER schemas into this model is discussed in [6]. Kripperndorf and Song [12] address the problem of mapping a star schema into an ER schema, and Song et al. [18] analyze many-to-many relationships between fact and dimension tables. An extension of the ER model for capturing the dimensional paradigm is presented in [17].

Several methodologies for data warehouse design have also been proposed. The design of data warehouses according to the dimensional fact model is described in [7]. Object-oriented-based methodologies are presented by Giovinazzo [4] and by Trujillo et al. [19]. In [8], the authors propose an approach to data warehouse design that resembles the traditional database design process. Like in [1] and [2], they advocate that at least a conceptual or logical modeling activity should precede the actual implementation of a data warehouse.

In [10] and [11], Kimball et al. provide a detailed discussion of the dimensional model concepts and use. In [10], a complete framework for applying the dimensional modeling approach, called Business Dimensional Lifecycle, is described. From the authors' perspective, requirements definition constitutes a fundamental step as it has a major impact on the whole proposed lifecycle. They suggest interviews and facilitated sessions as the basic techniques for this step. The difficulty of identifying dimensions is also addressed by Kimball in [9], and considered to be mostly a data warehouse designer's decision and responsibility.

To face the problem of decision support requirements being indeterminate and unstable, Lambert [13] presents suggestions on how to get users more involved, but yet with no direct support for this interaction. Firestone [3] proposes a dimensional object modeling approach, which strongly relies on information system use cases tightly coupled with strategic goals and objectives. According to the author, "the use case provides a view of the system from the view point of its users and their business purposes, … and it follows that the whole system architecture will be determined by the use cases and ultimately by the users who specified them."

Although many of these works consider and even emphasize the user involvement on early stages of requirements analyzes, little effort has been put on developing tools to directly support this task. Many commercial solutions for software components of data warehouse environments address the definition and specification of dimensions, fact variables and hierarchies. Modeling tools like Computer Associates's AllFusion ERwin Data Modeler, Oracle's CASE 2000 and MicroStrategy's Architect have comprehensive features to represent and describe star and snowflake schemas, but they do not include facilities for the participation of end-users in early stages of development. The same happens with ETL systems and OLAP tools, like SAS Administrator, Microsoft DTS and Microsoft OLAP Services, which contemplate later stages of data warehouse development. Thus, MD2 may be seen as an attempt towards supporting the user involvement in the development of data warehouse applications.

## 3 Overview of MD2

In this section, we present an overview of the MD2 tool. MD2 has been designed to work along with other tools that comprise a typical data warehouse architecture. Figure 1 depicts the MD2 environment, illustrating how its main components and the actors involved are related to each other. This environment is briefly discussed next.

As we can see from Figure 1, we distinguish three actors in this scenario. The *information analyst* is responsible for the development of the data warehouse application. He is the one who inputs the MD2 data repository with a preliminary definition of the dimensions, hierarchies, facts, and attributes identified for the application. We notice that sometimes this role might be played by the data warehouse designer too. The *expert user* is an application specialist that is responsible for structuring the contents of the data repository according to a hierarchy of application related themes. In doing so, he also helps to identify new elements (dimensions, hierarchies, facts, and attributes) that are relevant to the application. He is also responsible for later analyzing the data warehouse contents (using some specific data access tool) in order to create and generate application reports that will be stored in the Web server to be made available through a Web interface. This interface is also generated by the expert user based on the thematic structures defined in the MD2 data repository. When performing these tasks, both the information analyst and the expert user may be assisted by a number of specific reports generated by MD2. Finally, the *application user* (or *end-user*) might be any person that, through the Internet or a company Intranet, accesses the application reports using a Web browser. Application users play an important role in this environment because they help the expert user to identify the hierarchy of themes according to which the data repository is structured.

A more detailed discussion of how the MD2 tool is used to assist the development of data warehouse applications is postponed to Section 4. In what follows, we describe the MD2 data repository and present an overview of the MD2 user interface.

## 3.1    The Data Repository

The MD2 data repository is used to hierarchically represent the application data in terms of thematic structures. Each thematic structure corresponds to a *domain* (main subject) that is hierarchically composed of *structured subjects*, optional *composed facts*, and *items.* These components are called *structural elements*. In addition to structural elements, *independent elements* are used to compose or characterize other elements. The independent elements represented in the data repository are: *subjects*, *facts*, *attributes*, *dimensions*, *dimension hierarchies, reports*, and *report types.* Except for subjects, reports and report types, these independent elements correspond to the usual dimensional modeling primitives [10,11]. Details of how such thematic structures and the corresponding elements are created and manipulated in the data repository are given next.
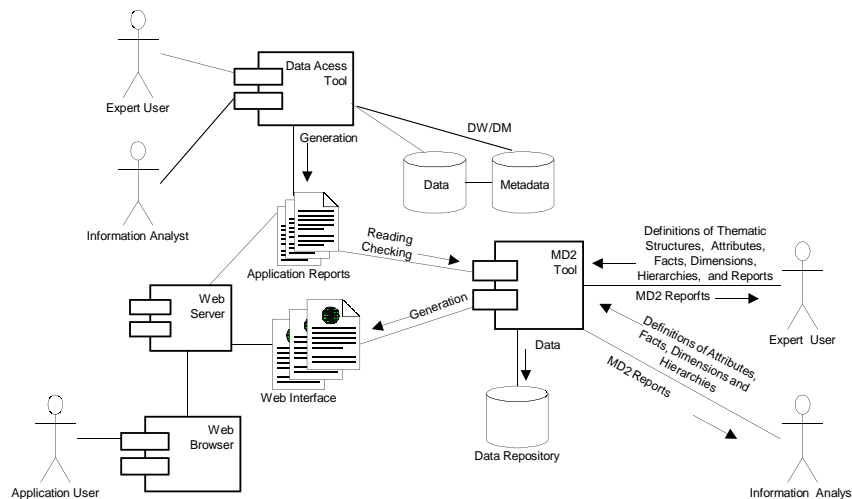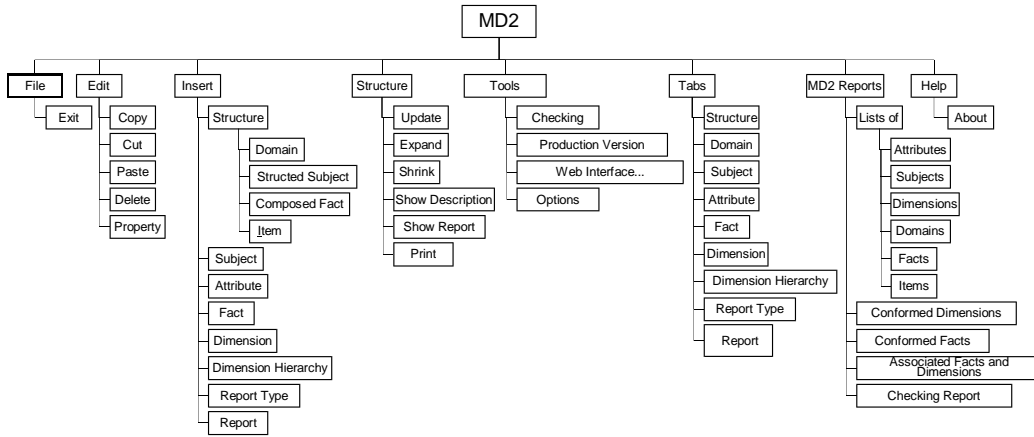


Figure 1: The MD2 environment.

Figure 2: MD2 function menu diagram.

## 3.2 The User Interface

MD2 provides a Windows-like graphical interface for user interaction. This interface has been implemented using standard Visual Basic classes [14] and includes function bars, pull-down menus, and specific function buttons. Figure 2 presents a diagram describing the tool's function menu. The upper level functions (File, Edit, Insert, Structure, Tools, Tabs, Reports, and Help) correspond to general functions that are included in the main screen bar. Clicking on any of these functions pops-up a pull-down menu for selecting a specific function.

The thematic structures and the elements described in the data repository (domains, subjects, facts, attributes, dimensions, dimension hierarchies, reports, and report types) are created and further manipulated by selecting specific tabs. For instance, by selecting the Estrutura (Structure) tab (see Figure 3), the user can create or modify a thematic structure by inserting, updating, and deleting its structural elements (domains, structured subjects, composed facts, and items). In addition, this function also allows the user to visualize some of the independent elements that might be associated with the thematic structure (reports associated with an item, dimensions associated with a report, dimension hierarchies associated with a dimension, and attributes associated with a dimension hierarchy).

Analogously, by selecting any other tab the user has access to the corresponding element screen. For example, Figure 4 shows the fact screen that is selected by clicking on the Fato (Fact) tab (see Figure 3). Notice that when a fact in the list box is selected, the fact properties (name, description, type, default aggregation rule, etc.) are presented in the text boxes below. The buttons Inserir (Insert), Alterar Propriedades (Update Properties), and Excluir (Delete) allow the user to, respectively, insert a new fact, update the properties of an existing fact, and delete an existing fact.
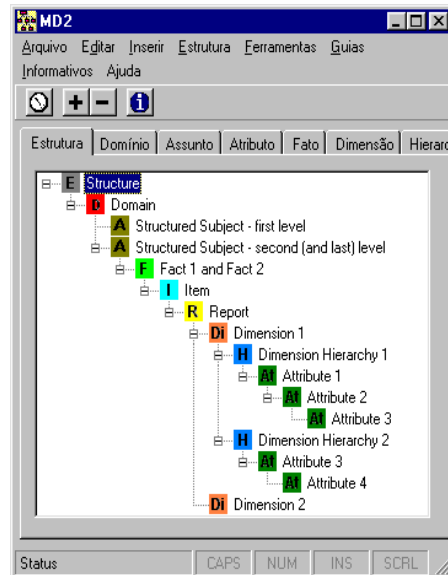


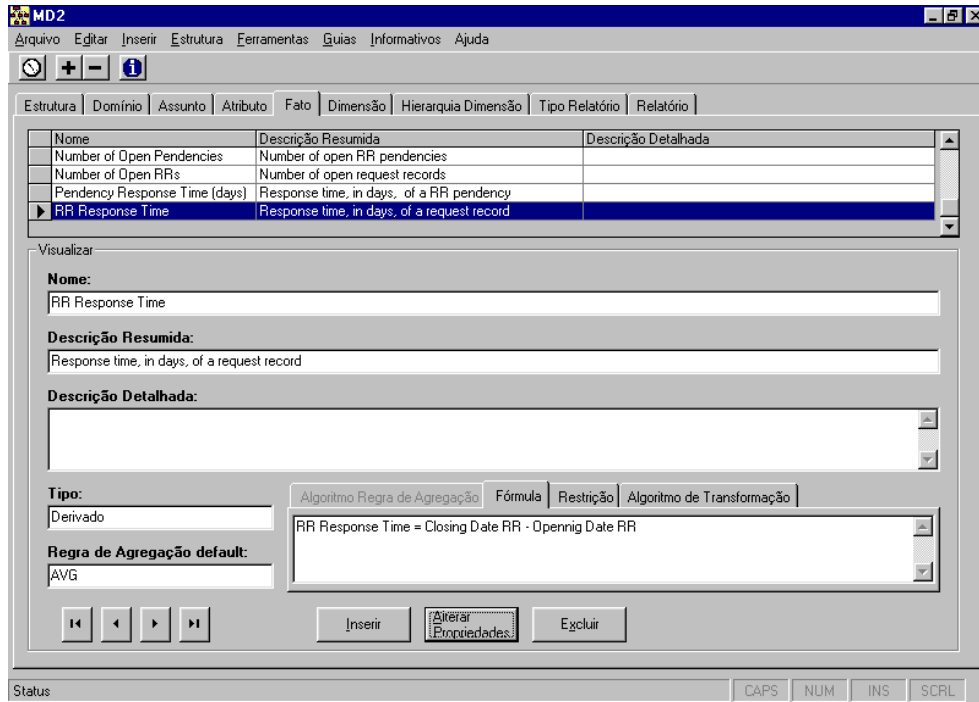Figure 3: Elements of a thematic structure.

Figure 4: Fact screen

The menus **Ferramentas** (Tools) and **Informativos** (MD2 Reports) provide additional facilities for supporting the work of the information analyst and of the expert users when modeling the application. These facilities include functions for checking the consistency of thematic structures and elements represented in the data repository, as well as for producing analytic reports to support some of the steps in the data warehouse developing process and generating a Web interface to publish the application reports that will be extracted from the data warehouse. Figure 5 presents the Web interface generated for the structure shown in Figure 3. The interface includes a navigation area and a presentation area. The navigation area displays the thematic structure allowing the user to navigate through its constituent elements. When an item is selected, the associated report is displayed in the presentation area. The command bar includes buttons for expanding and shrinking the navigation structure, activating the floating menu, and selecting the navigation mode (by item or dimension).
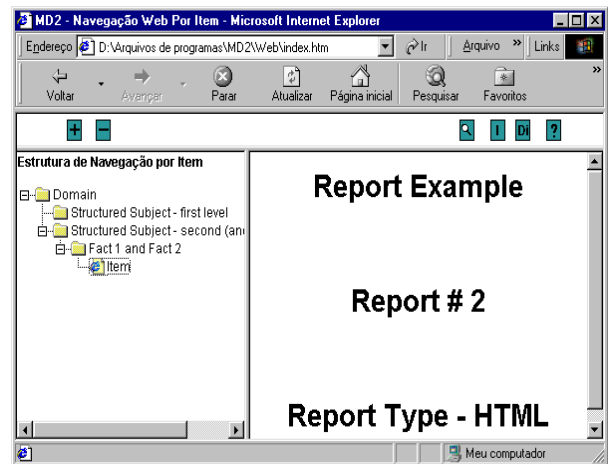


Figure 5: The Web interface

## 4 MD2 and the Business Dimensional Lifecycle

Developing a data warehouse application is a complex and time-consuming process that involves several steps. The MD2 tool aims at supporting some of these steps. To illustrate this, in this section we briefly discuss how it can be used in the context of

the framework provided by the Business Dimensional Lifecycle. Figure 6 describes the basic steps of this methodology as proposed in [10]. The shadow boxes correspond to the methodology steps in which MD2 can be applied.

In the Business Requirement Definition step, the tool can be used by the expert user to express, in terms of the dimensional model, his data analysis needs, facilitating a better understanding of the facts, dimensions, hierarchies, and attributes involved in the application. It also allows other thematic structures to be queried, revealing new elements that might be useful for data analysis.

The Dimensional Modeling step can be considerably facilitated by the contents of the MD2 data repository input in the previous step. From the thematic structures described in the data repository, the information analyst can identify possible data marts. As a first assumption, we assume that the last level of each subject corresponds to a data mart. It is also possible to identify associated facts and dimensions, descendents of a structured subject, and dimensions used to aggregate a fact. Table 1 lists the MD2 reports that can help execute the activities of this step.

In the Data Staging Design and Development step, some MD2 reports (e.g., List of Facts, List of Attributes, and List of Dimensions) can be used to help the information analyst create the data staging area metadata, especially the definition of attribute sources and formulas,

algorithms and transformations rules for attributes and facts.

MD2 can also assist some activities during the End-User Application Specification and Development steps. Particularly, the Web interface generated based on the thematic structures described in the data repository can be seen as an end-user application. Information in the List of Reports, especially on dimensions and related hierarchy attributes, report types, constraints and parameters, can help the specification of new application reports. In addition, other MD2 reports listing elements such subjects, dimensions, attributes and facts can help populate metadata for data access tools.

Finally, in the Maintenance and Growth step, the expert user can input in the data repository new elements that represent new requirements identified for the application.

## 5 Case Study

In this section, we briefly describe a case study in which the MD2 tool was used to design and develop a small data warehouse application. This application is a data mart prototype developed for BHTRANS, the Transport and Traffic Authority of Belo Horizonte, the third largest city in Brazil. The use of MD2 in this application followed the steps of the Business Dimensional Lifecycle, as discussed in the previous section.
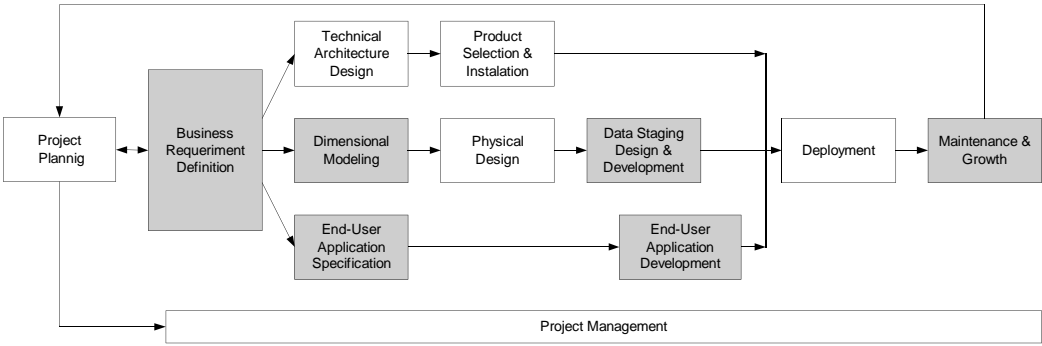


Figure 6: Steps of the Business Dimensional Lifecycle.

| Activities in the Dimensional Modeling step of the Business Dimensional Lifecycle | MD2 Reports | | | | | | |
|---|---|---|---|---|---|---|---|
| | Dimensions | Facts | Reports | Conformed Dimensions | Associated Facts and Dimensions | Subjects and Domains | Attributes |
| Building the matrix of data marts and dimensions | | | | | X | | |
| Defining the fact table granularity and selection of the dimensions according to the defined granularity | X | | | | X | | |
| Selecting facts | | X | | | X | | |
| Drawing the fact table diagram | | | | | X | | |
| Detailing the fact tables | | X | | | | | X |
| Detailing the dimension tables | X | | | | | | X |
| Mapping data from source to target | X | X | | | | | X |

Table 1: MD2 reports and activities in the Dimensional Modeling step.

This data mart corresponds to the Request Record, a subject related to the domain Citizen Enquiry. The main data source for this application is the system that records the requests, suggestions and complains from the city's transport system users. These requests, suggestions and complains, generally known as Request Records (RRs), are recorded from distinct sources (e.g., phone calls, letters, electronic forms, etc.) and stored in an Oracle database. An RR is classified according to its subject and type of occurrence. In addition, every RR is associated with an identified user. An RR can be sent to distinct administrative sectors and while it has not yet received an answer it is considered as pendent. At any time, it is possible to follow the status of an RR. When an RR is closed, its final status is recorded in the database and an answer is sent to the user who made the request, suggestion or complain from which it has been issued.

Several users from different administrative sectors of BHTRANS participated in the development of this application. After accomplishing the Project Planning step (see Figure 6), MD2 was used next in the Business Requirement Definition step for expressing, in terms of the dimensional model, the application data analysis needs. The main dimensional model elements were identified by analyzing the database entity-relationship schema, and the attributes, facts, dimensions and hierarchies identified were input in the MD2 data repository by the information analyst. Figure 7 illustrates a form used to insert the fact RR Response Time (days). This fact has AVG (average) as its default aggregate rule. In addition, since it is only considered for RRs that have been closed, a constraint Status RR = 'CLOSED' is imposed on it. This fact has also been classified as Derivado (Derived) because it is the result of a formula.

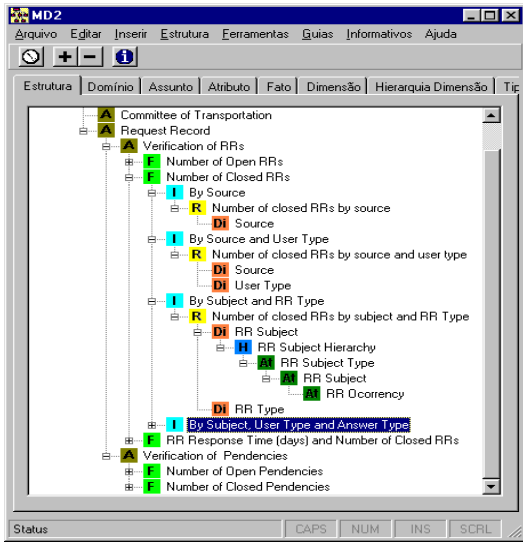

Figure 7: Insert Fact Form.

Figura 8: An application thematic structure.

After the input of the elements in the data repository, the application thematic structure was defined by an expert user. Figure 8 presents part of this structure that corresponds to the domain Citizen Enquiry and to the first level subject Request Record, which has as its descendants the subjects Verification of RRs and Verification of Pendencies. Figure 8 also illustrates how the composed facts (nodes F) are defined in terms of items (nodes I), which reports (nodes R) are related to each item, and the dimensions (nodes DI), hierarchies (nodes H) and attributes (nodes At) that compose each item. For instance, the subject Verification of RRs includes three of

composed facts: Number of Open RRs, Number Closed RRs, and RR Response Time (days) and Number of Closed RRs.

Then, the Dimensional Modeling step was carried out. In addition to information gathered from interviews in the previous step, several MD2 reports were used to help build the dimensional schema. Particularly, the report Associated Facts and Dimensions was used to validate the schema consistency and the report Conformed Dimensions was used to check the dimensions related to the two fact tables. Figure 9 depicts the application dimensional schema. The use of MD2 reports in the Data Staging Design and Developing step was not significant in this application.

In the End-User Application Specification and Development steps, several MD2 reports were used to help input metadata for the Oracle Discoverer [16] data access tool. For this task, the List of Reports was especially important to identify which dimensions and respective hierarchies should be extracted from the data warehouse to generate the application reports. The reports generated by Oracle Discoverer were then stored in the Web server for later access by end-users. To make possible the access to these reports, the Web interface was generated by an expert user. Figure 10 illustrates an application report generated through Oracle Discoverer and displayed by the Web interface.
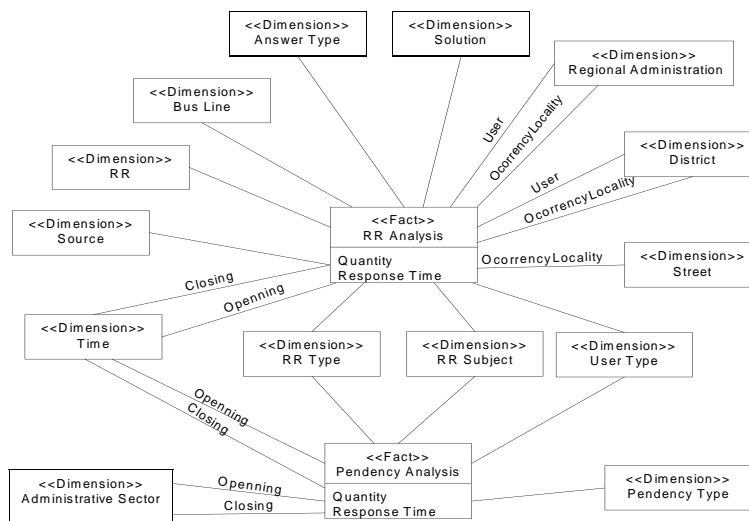


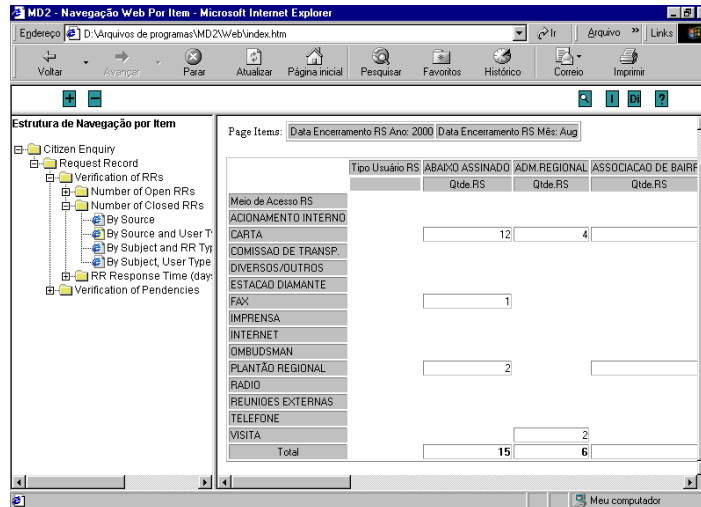Figura 9:  Dimensional schema of the data mart Request Record.

Figura 10: HTML report generated through Oracle Discoverer and displayed by the Web interface.

# 6 Conclusions

Any successful data warehouse initiative is first and foremost dependent on its business users. They must be engaged in the project from the very beginning as well as their excitement and motivation maintained along the whole development process. Most of all, they play a fundamental role during the requirements analysis and must be actively involved on the elucidation of relevant facts and perspectives.

We have presented in this paper MD2, a tool based on the dimensional modeling approach that assists the users' participation in the development of data warehouse applications. The tool is supported by a data repository and provides a user friendly graphical interface that helps to gather metadata for describing the dimensional schema as well as to specify reports that can be generated from the data warehouse and visualized through a Web interface. In addition, the tool includes a number of reports that help users perform their activities throughout the process of developing a data warehouse application.

To illustrate the use of MD2, we have discussed a case study within the framework of the Business Dimensional Lifecycle [10]. This case study has been based on a data mart prototype developed for BHTRANS, the Transport and Traffic Authority of Belo Horizonte, and emphasizes the user participation in some of the key steps of this methodology. However, it should be pointed out that the use of MD2 is not centered on the Business Dimensional Lifecycle steps and, therefore, the tool can be used with any other data warehouse design methodology based on the dimensional modeling approach, such as [8] and [15]. Particularly, MD2 might be used to capture the dimensional modeling elements derived from an entity-relationship schema according to the method proposed in [15].

As the MD2 tool is conceived to complement and interoperate with other tools usually found on data warehouse environments, the adherence to metadata standards [4] currently being discussed by OMG and other groups is a fundamental requirement. We believe the MD2 data repository metamodel can be easily adapted in the future to any unified standard that might result from current metadata standard initiatives.

## References

[1]   Agrawal, R., Guppta, A., and Sarawagi, S. Modelling Multidimensional Databases. *Proceedings of the Thirteenth International Conference on Data Engineering*, Birmingham, UK, 1997, pp. 232-243.

[2] Cabibbo, L., and Torlone, R. A Logical Approach to Multidimensional Databases. *Proc. of the 6$^{th}$ Int'l Conference on Extended Database Technology*, Valencia, Spain, 1998, pp. 187-197.

[3] Firestone, J.M. *Dimensional Object Modeling*. Executive Information Systems, Inc., White Paper n. 7, April 1998 (available at http://www.dkms.com/ DOM.htm).

[4] Giovinazzo, W. A. *Object-Oriented Data Warehouse Design.* Prentice Hall, New Jersey, NJ, 2000.

[5] Golfarelli, M., Maio, D., and Rizzi, S. The dimensional fact model: a conceptual model for data warehouses. *International Journal of Cooperative Information Systems 7*, 2-3 (1998), 215-247.

[6] Golfarelli, M., Maio, D., and Rizzi, S. Conceptual Design of Data Warehouses from E/R Schemas. *Proc. of the 31$^{st}$ Hawaii Int'l Conference on System Sciences, Vol. VII*, Kona, Hawaii, 1998, pp. 334-343.

[7] Golfarelli, M., and Rizzi, S. Designing data warehouses: key steps and crucial issues. *Journal of Computer Science and Information Management 2*, 3 (1999).

[8] Hüsemann, B., Lechtenbörger, J., and Vossen, G. Conceptual Data Warehouse Design. *Proc. of the Int'l Workshop on Design and Management of Data Warehouses*, Stockholm, Sweden, 2000, pp. 6.1-6.11.

[9] Kimball, R. Mystery Dimensions. *Intelligent Enterprise Magazine 3*, 5 (March 2000).

[10] Kimball, R., Reeves, L., Ross, M., and Thomthwaite, W. *The Data Warehouse Lifecycle Toolkit: Tools and Techniques for Designing, Developing and Deploying Data Warehouses.* John Wiley & Sons, New York, 1998.

[11] Kimball, R. A Dimensional Modeling Manifesto. *DBMS 10*, 9 (August 1997).

[12] Krippendorf, M., and Song, I.-Y. The Translation of Star Schema into Entity Relationship Diagrams. *Proc. of the Eighth Int'l Workshop on Database and Expert Systems Applications, DEXA'97*, Toulouse, France, 1997, pp. 390-395.

[13] Lambert, B. Break Old Habits To Define Data Warehousing Requirements. *Data Management Review* (December 1995).

[14] Microsoft. *Visual Basic Developer's Guide.* 1999.

[15] Moody, L.D., and Kortink, M.A.R. From Enterprise Models to Dimensional Models: A Methodology for Data Warehouses and Data Mart Design. *Proc. of the Int'l Workshop on Design and Management of Data Warehouses*, Stockholm, Sweden, 2000, pp. 5.1-5.12.

[16] Oracle. *Oracle Discoverer 3.1 Release 3.1: Administration Guide*. 1998.

[17] Sapia, C., Blaschka, M., Höfling, G., and Dinter, B. Extending the E/R Model for the Multidimensional Paradigm. *Proc. of the Int'l Workshop on Data Warehousing and Data Mining*, Singapore, 1998, pp. 105-116.

[18] Song, I.-Y., Rowen, W., Medsker, C., and Ewen, E. An Analysis of Many-to-Many Relationships Between Fact and Dimension Tables in Dimension Modeling. *Proc. of the Int'l Workshop on Design and Management of Data Warehouses*, Interlaken, Switzerland, 2001, pp. 6.1-6.13.

[19] Trujillo, J., Palomar, M., Gómez, J., and Song, I.-Y. Designing Data Warehouses with OO Conceptual Models. *IEEE Computer 34*, 12 (2001), 66-75.

[20] Tryfona, N., Busborg, F., and Christiansen, J. starER: A Conceptual Model for Data Warehousing Design. *Proc. of the ACM Second Int'l Workshop on Data Warehousing and OLAP*, Kansas City, MI, 1999, pp. 3-8.