

Alternative Development for Data Migration Using Dynamic Query Generation

Johan Alfredo Romero-Ramírez¹, Carlos Enrique Montenegro-Marín², Vicente García-Díaz³,
Juan Manuel Cueva Lovelle⁴

^{1,2} Engineering Faculty, Universidad Distrital Francisco José de Caldas, Colombia,
^{3,4} Computer Science Department, University of Oviedo, Spain

Abstract – This article presents an ETL (Extract, Transform, Load) prototype called Valery as alternative approach to migration process which includes a compiler for dynamic generation of SQL queries. Its main features involve: SQL dynamic generation, set of configuration commands and environment for file uploading. The tests use the Northwind academic database and an individual environment. The model implementation uses flat files and SQL as query language. Finally, there is an analysis of the results obtained.

Keywords – Data migration, Data transformation, ETL, schema mapping.

I. INTRODUCTION

John Morris (2001) [1] defines the data migration as the selection, preparation, extraction, transformation and movement of the suitable data with the right quality, to the right place and with the dismantling of the legacy data sources.

Some of these projects may have the scenario where the information sources involve flat files [7] and the target system is a relational database and SQL-based [9].

This specific domain presents the problem that usually structures and formats of the legacy data differ or are incompatible with the relational design of the target system and in some cases unexpected changes are raised in these structures.

Therefore, the process must resolve two common and complex tasks in the migration process [8]: (i) to associate the different scheme mapping related elements between the source and the destination (*mapping*), (ii) to react effectively to the specification schema changes (*schema changes*).

Figure 1 shows an example of different sources. The scheme of the legacy system presents a different format from the target system related tables and the arrows represent the correspondence between the source-target scheme fields. It is possible to see the differences between the corresponding formats.

This paper presents the features of a solution called *Valery* as an alternative approach to data migration for heterogeneous data sources where the source implies flat files. This is based on the hypothesis that it is possible to create a model that exploits the advantages of standardisation and functionality existing by SQL language (Structured Query Language) [10] and the RDMB (Relational Database Management System) [11] onto a model that addresses the problem of mapping exposed and delegates the factors of performance to the infrastructure hardware and legacy data size.

II. FEATURES

Valery approach has the following characteristics:

Conversion of the flat file into a temporary table in the system target: Flat files are converted to a temporary table in the database manager. The advantage of this approach enables subsequent handling by the SQL statements in a transparent manner between the schema source and destination.

Specification of generic cases for the mapping between the schema source and destination: Mapping patterns are identified to encapsulate the correspondences between the data sources and target tables. This allows adapting the model to a wide range of migration and integration cases which comply with the generic features and assure independence in the relational structure.

Specification of SQL statement associated with each case of migration: Different SQL queries needed to solve each generic case for migration and integration are identified.

Generation of dynamic queries: Represents the components associated with the generation of dynamic SQL queries. Query Builder creates dynamically the sentences due to generic pattern recognition and consistent configuration of the destination relational model. For this purpose, *Valery* has a

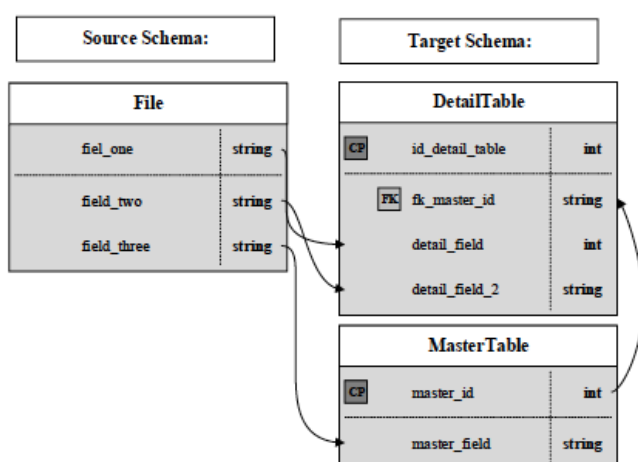


Fig. 1. Schema representation with a non-linear connection.

command console that receives the parameters associated with a specific migration project.

Evaluation with specific data: Implementation was carried out using the *NorthWind* model and own model. The results indicate a satisfactory performance in quality and runtime.

III. VALERY MIGRATION PROCESS

Mapping specification: The first step is to configure the source with the relational database format. This is an interface that allows capturing such a configuration in explicit defined commands.

Load source files: In this step, loading flat files for validation and migration of the flat file in a new temporary table within the database target system is made for subsequent use in the query execution.

Creation of queries: The system generates the specific SQL queries according to the settings made.

Execution of SQL queries: Running SQL queries created and displaying the report of results thrown by the database manager.



Fig. 2. Valery process.

IV. STATE OF THE ART

Different studies have been performed for the automatic execution of mapping schemas. To determine the correlation of patterns these studies propose techniques that exploit different types of information, such as definitions, data types, structures of schemas and instances of data [2].

As reference studies in the development of some components of the proposed model we have presented *QuickMig* tools [2] and *HumMer* (Humboldt Merger) [3]. *QuickMig* is a new semi-automatic approach for the identification of semantic correspondences between the sources of information schema elements. *QuickMig* presents a series of new techniques that exploit example cases, domain ontologies and the re-use of existing allocations to detect not only the correlation of elements but also their expressions of mapping. *QuickMig* includes new mechanisms to effectively incorporate the knowledge of the domain's users in the comparison process. The results give a full evaluation using diagrams and actual data [2]. On the other hand, *Hummer* is a tool *ad-hoc* that allows for the fusion automatic of data using extensions SQL for heterogeneous schemes with data duplicate or in conflict.

V. OVERVIEW

This section provides a review of *Valery* approach to the correlation of schemes.

A. Considerations

Valery implementation requires knowledge of the domain of the migration process. As also outlined in [2], migration processes need knowledge about the origin and the target system. Unfortunately, such knowledge may not be always available in one place. Knowledge about the target system may be available in the provider, while only customers can provide a detailed origin systems knowledge.

B. Scope

The solution is applicable to projects that could get information bequeathed as flat files and where there is SQL-based database server (RDMB).

Valery offers a generic case of migration domain. These cases include data with foreign keys, different types of data, updating data processes, independence from flat file formats and independence of the database destination complex.

C. Cases

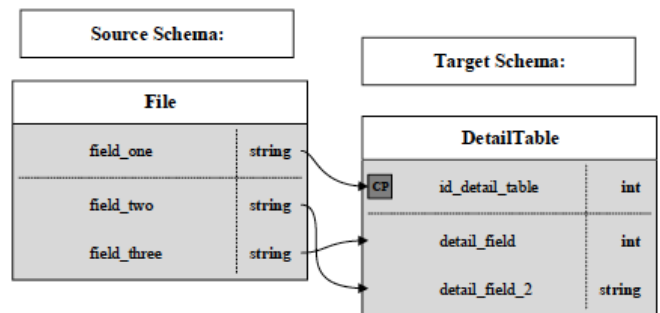


Fig. 3. Schema representation with a linear correspondence.

Linear correspondence of the flat file and database structures that do not include relational data. In Fig. 3 it is possible to see an example of considered cases. The associated fields have a linear relationship with the destination schema and show differences in nomenclature and data type.

There is a linear correlation between the flat file and the structures of the database due to referential integrity presenting tables: the original file must populate all tables in a single transaction (see Fig. 1). Data are generated massively associating foreign keys dynamically. In addition, there are differences in nomenclature and data type.

D. Algorithmic Theoretical Foundation

Valery provides a set of rules for the expression of associated semantic generation of SQL queries and where should be *correlated* data from the data source and the database file destination. The system makes reading of the commands and generates the corresponding sentence type *INSERT-SELECT-UPDATE* by a generator queries (*Query Engine*). Script specification depends on the complexity of the destination database and of the different levels of the data base relational hierarchy.

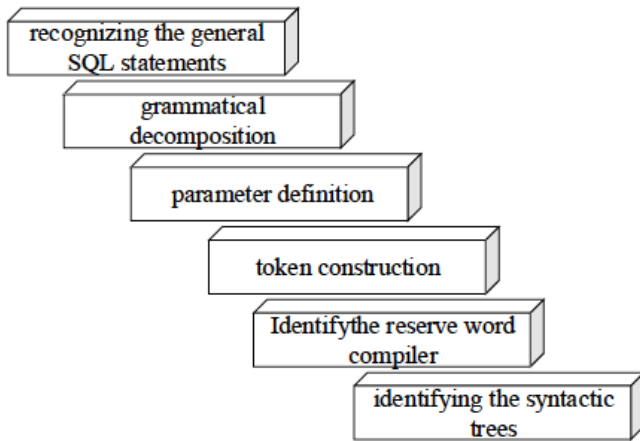


Fig. 4. Valery methodology.

The methodology for the generation of SQL queries collects (see Fig. 4): (i) recognition of the SQL statements associated with each case, (ii) decomposition grammar of each sentence, (iii) definition of parameters, (iv) valid from the reading of the sentences *tokens* construction target, (v) creation of the keywords of the compiler according to the parameters, (vi) identification of syntactic trees. This allows for one theoretical basis for the modelling and design of the query engine that will take as input the configuration commands according to a specific migration project.

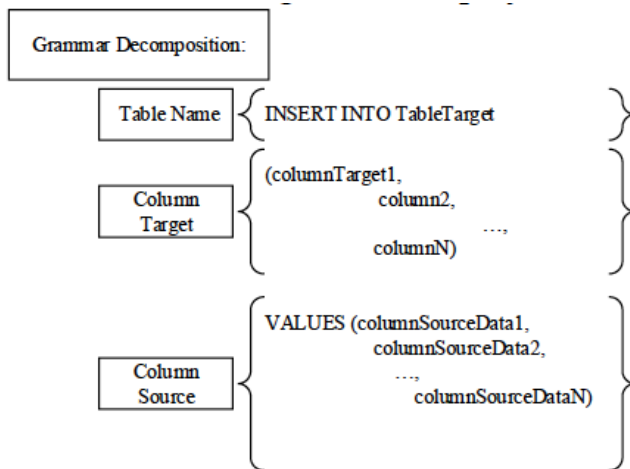


Fig. 5. Grammatical case of inclusion of one decomposition.

As an example, the grammatical breakdown of target SQL query associated with case (1) in Fig. 5 is displayed.

The model breaks down the previous sentence in the following grammatical fragments that symbolise the commands in the system for their respective configuration:

- Table name.
- Set of fields in the table destination.
- Set of fields in the table origin.

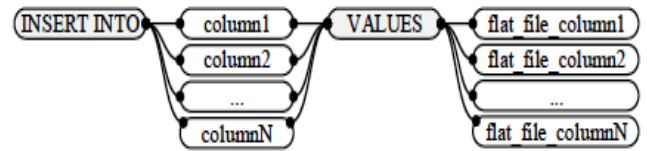


Fig. 6. Graph of an SQL statement.

Graph (Fig. 6) representation and representation in first order algebra allow obtaining the following relevant attributes for the *Query Builder* design.

The graph shows the following relevant properties for dynamic generation algorithm:

- The root node represents the keyword INSERT INTO.
- The next leaf node presents the schema columns.
- Equivalent to the keyword VALUES inner node is equal to the number of columns of the template.
- The length of the path from the root node to the leaf node is proportional to the number of columns.

On the other hand, a representation in first order algebra can be represented to a mathematical formalization of the objective statement:

$$TargetTable \leftarrow TargetTable \cup \{(value1, Value2, \dots, ValueN)\}$$

Due to the previous theoretical foundation, it is possible to represent the attributes required in script type attribute-value for the correct generation of case one target SQL statements. Below the commands are shown that are required by the *QE*:

Attribute (Command_TT): value (name of the destination table).

Attribute (Command_FS): value (field in the source table).

Attribute (Command_FT): value (field of the destination table).

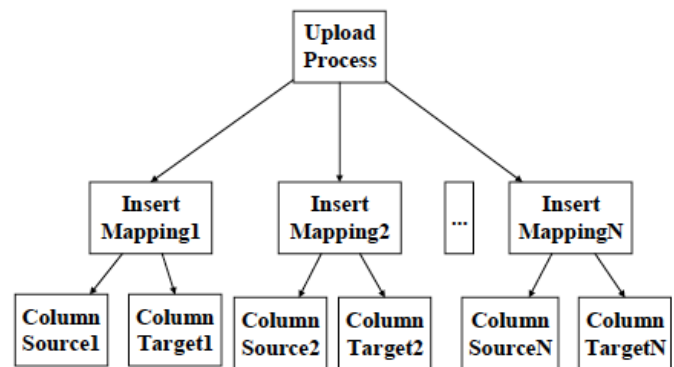


Fig. 7. Binary tree used by the Query Builder.

In addition, the *QE* algorithmic design implements a data structure of the syntactic tree such as a perfect binary tree where the hierarchy of commands depends on the foreign keys of tables. Figure 7 represents the syntactic tree for a specific case where the representation of commands is modelled.

It is noted that the route of the syntactic tree does not require rearrangement or searches to meet a predefined SQL

syntax [4], is only requirement the algorithm implemented functionality of insertion into the binary tree elements.

This tree insertion in the binary tree has an $O(\log n)$ complexity [5] and is recognised as a problem class P (treatable) [6].

VI. VALERY ARCHITECTURE

Valery components include: (i) **use of flat files**: flat files as legacy information due to the support that lies in the integration of various applications in these formats. (ii) **GUI**: this layer deploys all components related to the design, data capture, the configured schema validation, loading flat files, reports of the processes and operations carried out by the users to interact with the system. (iii) **Query Engine**: represents the components associated with the dynamic generation of SQL statements. To do this, follow the exposed methodology (Section V.D). The *QE* is responsible for compiling the predefined rules and generating SQL statements according to the request made by the user. This client/server architecture allows for *scalability* and *concurrency* of multiple migration processes. (iv) **Data access layer**: the components dedicated to establish the connection with the RDBMs, affect the target database according to each charging process and the results of the application.

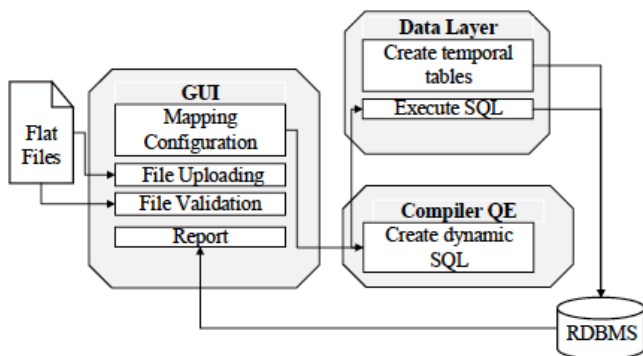


Fig. 8. Valery architecture model.

Valery follows the architecture presented in Fig. 8. The core (*QE*) is based on the grammatical breakdown of SQL statements that satisfy the cases of migration and integration. A set of commands and rules where the *QE* works as an interpreter of these commands is defined in the GUI.

VII. TESTING SCENARIO

Test database (*Northwind* and the own model) was used as target source, SQL Server as RDBMS and several flat files (*.csv) structurally separated with table target. With this testing scenario created the following accounts were verified:

1) Transformations of data complexity: For this purpose a set of migration process affecting different tables using as input files that have structural independence with the target database was defined.

2) Loading the flat file as an additional table in the database destination: this leads to a better performance of the SQL queries execution.

3) Tracking and audit: the statements to be executed in the respective RDBMS use all the functionality offered by the RDBMS for tracking and audit.

Technical hardware and software for testing: Intel® Core™ 2 Duo Processor, Windows 7 Home Premium (64-bit), 4 GB DDR2 800 MHz memory, 500GB HDD (5400 RPM, Serial ATA), Sql Server 2008®.

Migration *Northwind* processes involved 32 SQL queries generation while that for the own database processes involved 16 SQL queries generation. The data types of the formats included strings, dates and numeric values.

Example 1. INSERT queries generated for the Northwind database.

```

INSERT INTO Categories (
  CategoryName,Description,Picture)
SELECT cmpExcel5,cmpExcel8,cmpExcel11
FROM TablaExternaV10
  
```

Example 2. INSERT queries generated for the Northwind database.

```

INSERT INTO Products (
  SupplierID,ProductName,QuantityPerUnit,Un
  itPrice,UnitsInStock,UnitsOnOrder,ReorderLe
  vel,Discontinued)
SELECT
  (SELECT TOP 1 tabla0.SupplierID FROM Suppliers AS
  tabla0
  WHERE tabla0.ContactName = cmpExcel25),
  cmpExcel28,cmpExcel31,cmpExcel34,cmpExc
  el37,cmpExcel40,cmpExcel43,cmpExcel46
FROM TablaExternaV10
  
```

Example 3. UPDATE queries generated for the Northwind database.

```

UPDATE Orders
SET Freight = campoExcel81,
  ShipName = campoExcel84
FROM TablaExternaV10
WHERE OrderID =
  (SELECT TOP 1 tabla0.OrderID
  FROM Orders AS tabla0,
  Customers AS tabla1,
  Employees AS tabla2
  WHERE tabla0.ShipVia = cmbDinamico51
  AND tabla1.CustomerID = tabla0.CustomerID AND
  tabla1.ContactName = cmpExcel59
  AND tabla2.EmployeeID = tabla0.EmployeeID AND
  tabla2.FirstName = campoExcel66)
  
```

Table I and Table II show the results of the tests according to the model of testing used.

TABLE I
EVALUATION OF RESULTS OF THE QUERY EXECUTION

BD	Queries	Time (sg)
NorthWind	32	0.54
SAP	16	0.26

TABLE II
AVERAGE RESPONSE OF PROCESSES ACCORDING TO FILE SIZE

Size	Average: Upload and Execution
1 mb	12 sg
10 mb	249 sg
100 mb	3526 sg
1000 mb	17513 sg

In Table I, it is possible to observe performance in the generation of SQL queries. It obtains the expected result of the syntax (see examples 1, 2, 3) in the generation of queries that will be subsequently executed in the RDBMS.

On the other hand in Table II, it is possible to see that the scenarios where the files have a weight of larger system present an unfavorable performance. This is due to the conversion of the flat file into the temporary table that is proportional to the size of the file and the query execution that is also proportional to size of the temporary table.

VIII. CONCLUSION AND FUTURE STUDIES

The approach presented by Valery was experimentally evaluated using data in a test environment. It is quite effective for creating well-formed SQL queries according to the defined specific domain. Moreover, proportional migration performance depends of the size file and the database server hardware. Therefore, performance is an independent factor to the generation of SQL queries.

Our approach using generic cases allowed for a coverage of different sorts of migration requirements in two testing databases using insertion and update.

Mapping and data transformation conditions could solve specific tests. The approach allowed encapsulating the problem of mapping settings and delegating the performance issue to the hardware infrastructure and file sizes corroborating the initial hypothesis.

The future studies will be connected with similar problems arising in the model transformation tasks [12]–[16].

REFERENCES

- [1] J. Morris, *Practical Data Migration*, BCS. 2012, pp. 7, 8, 9, 10, 77, 164.
- [2] C. Drumm, M. Schmitt, H. H. Do and E. Rahm, "QuickMig. Automatic Schema Matching for Data Migration Projects," in *Proc. of the sixteenth ACM conf. on Conf. on information and knowledge management*, CIKM'07, ACM New York, NY, USA, 2007, pp. 107–116. <http://dx.doi.org/10.1145/1321440.1321458>
- [3] A. Bilke, J. Bleiholder, C. Böhm, K. Draba, F. Naumann, M. Weis, "Automatic Data Fusion with HumMer," in *Proc. of the 31st int. conf. on Very large data bases*, VLDB '05, pp. 1251–1254, 2005.
- [4] V. Ebaï, *What Is Sql?: Fundamentals of Sql, T-Sql, PL/Sql and Datawarehousing*, 2012, pp. XI.
- [5] K. Loudon, *Mastering Algorithms with C*, O Reilly Media, Inc. 2009, p. 206.
- [6] A. D. Munoz, *Metaheuristics*. Ed. Dykinson, 2007, p. 12.
- [7] B. R. Ullrey, *Implementing a Data Warehouse: A Methodology that Worked*. AuthorHouse, 2007, pp. 93–94.
- [8] Z. Bellahsene, A. Bonifati, E. Rahm, *Schema Matching and Mapping*. Springer Science & Business Media. 2011, pp. 152, 153.

- [9] A. D. Ionita, *Migrating Legacy Applications: Challenges in Service Oriented Architecture and Cloud Computing Environments*. 2012, pp. 210.
- [10] S. Kedar, *Database Management System*, Technical Publications. 2009, p. 42.
- [11] S. Sumathi, S. Esakkirajan, *Fundamentals of Relational Database Management Systems*. Springer, 2007, p. 26.
- [12] U. Donins, J. Osis, A. Slihte, E. Asnina and B. Gulbis, "Towards the Refinement of Topological Class Diagram as a Platform Independent Model," in J. Osis, O. Nikiforova (Eds.). *Model-Driven Architecture and Modeling-Driven Software Development*. ENASE 2011, 3rd Whs. MDA&MDS, SciTePress, Portugal, 2011, pp. 79–88.
- [13] J. Osis, E. Asnina, A. Grave, "Computation Independent Representation of the Problem Domain in MDA," *e-Informatica Software Engineering J.*, vol. 2, issue 1, 2008, pp. 29–46.
- [14] J. Osis, U. Donins, "Formalization of the UML Class Diagrams," *Evaluation of Novel Approaches to Software Engineering*. Springer-Verlag, Berlin Heidelberg, New York, 2010, pp. 180–192. http://dx.doi.org/10.1007/978-3-642-14819-4_13
- [15] J. Osis, A. Slihte, "Transforming Textual Use Cases to a Computation Independent Model," in J. Osis, O. Nikiforova (Eds.). *Model-Driven Architecture and Modeling Theory-Driven Development*, ENASE 2010, 2nd MDA&MTDD Whs., SciTePress, Portugal, 2010, pp. 33–42.
- [16] A. Slihte, J. Osis and U. Donins, "Knowledge Integration for Domain Modeling," in J. Osis, O. Nikiforova (Eds.). *Model-Driven Architecture and Modeling-Driven Software Development*. ENASE 2011, 3rd Whs. MDA&MDS, SciTePress, Portugal, 2011, pp. 46–56.

Johan Alfredo Romero Ramírez is a System Engineer of the Engineering Faculty of the Universidad Distrital Francisco José de Caldas. His research interests include domain-specific languages and web engineering. Contact address: Engineering Faculty of the Universidad Distrital Francisco José de Caldas, Carrera 7 N° 40–53 Piso 5, Bogotá (Cundinamarca, Colombia). E-mail: jaromerora@gmail.com

Carlos Enrique Montenegro-Marin is an Associate Professor at the Engineering Faculty of the Universidad Distrital Francisco José de Caldas. He has a PhD from the University of Oviedo in computer engineering. His research interests include model-driven engineering, domain-specific languages, technology for learning, cloud computing and programing languages. Contact address: Engineering Faculty of the Universidad Distrital Francisco José de Caldas, Carrera 7 N° 40–53 Piso 5, Bogotá (Cundinamarca, Colombia). E-mail: cemontenegrom@udistrital.edu.co

Vicente García-Díaz is an Associate Professor at the Computer Science Department of the University of Oviedo. He has a PhD from the University of Oviedo in computer engineering. His research interests include model-driven engineering, domain-specific languages, technology for learning and entertainment, project risk management, software development processes and practices. He is a Certified Associate in Project Management through the Project Management Institute. Contact address: Computer Science Department, University of Oviedo Edificio de la Facultad de Ciencias. C/ Calvo Sotelo s/n. 33007 Oviedo (Asturias, España). E-mail: garciavicente@uniovi.es

Juan Manuel Cueva Lovelle became a Mining Engineer at Oviedo Mining Engineers Technical School in 1983 (Oviedo University, Spain). He has a PhD from Madrid Polytechnic University, Spain (1990). From 1985 he has been a Professor in the languages and computers systems area at Oviedo University (Spain), and is an ACM and IEEE voting member. His research interests include object-oriented technology, language processors, human-computer interface, web engineering, modelling software with BPM, DSL and MDA. Contact address: Computer Science Department, University of Oviedo Edificio de la Facultad de Ciencias. C/ Calvo Sotelo s/n. 33007 Oviedo (Asturias, España). E-mail: cueva@uniovi.es