

Data Quality Improvement through Horizontal Fragmentation in Data Warehouses

Ljiljana Brkić, Igor Mekterović and Slaven Zakošek

Abstract— In this paper we propose the procedure for horizontal fragmentation of data warehouse tables and show how it affects data warehouse's data quality, namely, completeness and timeliness data quality dimensions. Horizontal fragmentation is suitable when data from multiple, more or less independent, organizational units are stored in the same fact table. In such cases, proposed procedure enables horizontal data separation based on the configured attributes thus isolating one institution's data quality problems from another. The proposed procedure relies on metadata and is generic and applicable to any data warehouse.

Keywords—Completeness, Data Warehouse, ETL, Horizontal fragmentation, Timeliness

I. INTRODUCTION

IT is common to describe and evaluate the quality of the data delivered to users in terms of quality attributes [1]. These attributes are called the data quality dimensions. Each dimension covers a specific aspect of quality. Data quality dimensions may relate to the value of the data or its purpose. Both, the data value quality and data schemes quality affect the quality of business processes. E.g., unnormalized schema of a table in relational data model results in redundancy and anomalies when inserting, updating and deleting data.

Data quality dimensions are usually defined from the qualitative aspect, referring to the general properties of the data values and data schemas without prescribing the way in which the quantitative values are assigned to dimensions. Despite the recognized importance of the quality of data schemes (both conceptual and logical), today, the focus is on the quality of data values [2]. For quantitative evaluation of quality dimensions it is

necessary to define metrics and measurement methods. The importance of this issue has been recognized and number of scientific papers (some of which are [1] [3] [4] [5]) deal with methods of measuring and quantifying the dimensions of data quality. There is no common position or accepted standard with regards to data quality dimensions classification or even data quality dimensions' definition. Thus, often for essentially the same data quality definition different terms are used, and vice versa, the same term has a different meaning in works by different authors.

Data quality in a data warehouse is affected not only by the quality of source data, but also by the quality of ETL process. For example, it is possible for an accurate and complete record from a data source to not be transferred to the data warehouse (or not timely transferred) due to faulty internal logic, inappropriate data warehouse refreshment policy or just plain errors in the ETL process. In other words, the data warehouse data quality dimensions depend on the quality of ETL processes.

The process of determining the quality of a data warehouse starts by assessing the quality of an individual attribute. The quality of a tuple is calculated by aggregating the quality of its attribute values. Analogously, the quality of table is determined via its tuples and, ultimately, with that information, the quality of the entire data warehouse is calculated.

In this paper we focus on two data quality dimensions: completeness and timeliness [6]. In order to achieve the improvement of those data quality dimensions, we propose and describe modifications to the data warehouse's ETL process.

II. HORIZONTAL FRAGMENTATION OF THE DATA WAREHOUSE TABLES

In this paper, we consider a generalized data warehousing architecture (Figure 1) consisting of three subsystems: data sources, staging area and data warehouse. Data warehouse systems that employ incremental [7] [8] or real-time [9] [10] refresh strategies must comprise a subsystem that is used to track changes in data sources in order to distinguish the new data needed to efficiently update the data warehouse. These subsystems are usually referred to as CDC (Change Data Capture) systems. Whether using CDC subsystem or not (for instance, on initial load), data is loaded into the

Manuscript received May 3, 2012.

This work was supported in part by the Ministry of Science, Education and Sports, Republic of Croatia under the project „Semantic integration of Heterogeneous Data Sources“

Lj. Brkić is with the Department of Applied Computing, Faculty of Electrical Engineering and Computing, University of Zagreb, Croatia (e-mail: ljiljana.brkic@fer.hr).

I. Mekterović is with the Department of Applied Computing, Faculty of Electrical Engineering and Computing, University of Zagreb, Croatia (e-mail: igor.mekterovic@fer.hr).

S. Zakošek is with the Department of Applied Computing, Faculty of Electrical Engineering and Computing, University of Zagreb, Croatia (e-mail: slaven.zakosek@fer.hr).

staging area, where it is transformed and subsequently loaded into the data warehouse in respect with the data warehouse schema. We assume that data warehouse employs relational or dimensional model, or both. In addition to metadata, the targeted data warehouse

contains transformed data in non-aggregated form (usually a relational database) and also more or less aggregated data (usually on OLAP server).

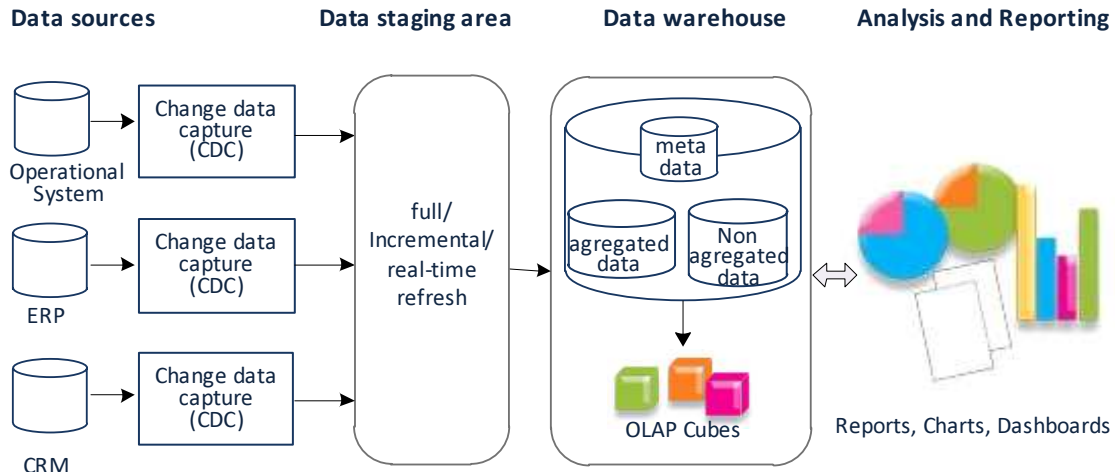


Figure 1 Data Warehouse System architecture

Data warehouses often integrate data of the same structure (schema) from several departments of a company, or a number of companies within a corporation, or a number of state institutions of a country, etc. An example of such environment is a data warehouse that integrates data from multiple higher education institutions such as the one described in [11]. The quality of the data belonging to a particular department, company or institution should primarily be a reflection of the effort that was put into the process. It is not unusual for one department to be more devoted to proper data administration than the other. Errors in the semantic integrity of the data belonging to one or several organizational units should not affect the others. The ETL process should be modeled and implemented with this in mind.

Due to the detected errors, the content of the fact table fExam has not been updated but was kept at the state from the previous update cycle (31.01.2012). Although there were no errors in the other institution’s data on 01.02.2012, errors in HEI 89’s data caused other institutions’ correct data not to be transferred. In other words, an institution is affected by the poor data quality belonging to other institution.

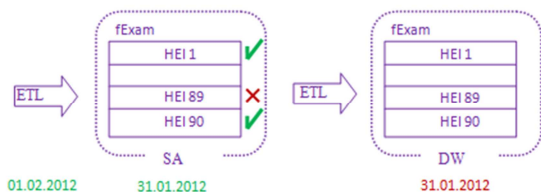


Figure 2 Errors in one fragment stop data of fragments

For instance, Figure 2 depicts update of fExam fact table of the data warehouse integrating data from 90 higher educational institutions (HEI 1 – HEI 90) in the Republic of Croatia [12] [13]. In this project, we had to employ an all-or-nothing refreshment policy; since it was of the utmost importance to have complete data (e.g. data warehouse is used to produce scholarship listings). Suppose that during the ETL process some integrity errors (e.g. foreign key violation) have been detected in records belonging to higher education institution HEI 89.

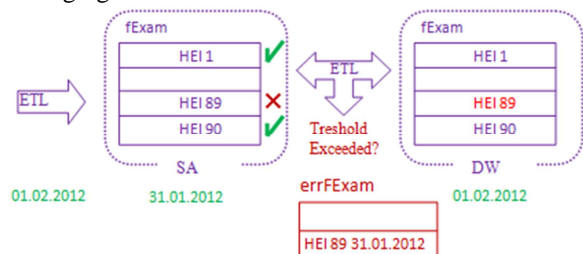


Figure 3 Errors in one fragment don't stop data of fragments

Figure 3 illustrates a better solution employing the following improvements:

- fExam table is logically divided to horizontal fragments, each fragment belonging to a particular higher education institution

- ETL process is modified as to be capable to find error free fragments in the staging area, and accordingly refresh corresponding fExam fragments. Fragments that exceed the error threshold (maximal allowed number of incorrect tuples) are not updated, i.e. remain as they were after on the last successful refreshment cycle.

This is much more robust and customizable approach. Each higher education institution is allowed to define the number of errors that it is willing to tolerate. If the defined threshold for a particular institution is exceeded,

the associated fragment of the fact table will not be updated. That fragment retains the state in which the number of errors was within the predefined threshold. Such strategy can lead to the state in which the different fragments of fact and dimension tables are updated in different update cycles. Tables become unions of approved (above threshold) fragments, with each fragment having its own timestamp. Diligent institutions (departments, etc.) are no more affected by others and are more likely to have up-to-date data in the data warehouse. On the other hand, if one wants to ignore errors in favor of timeliness, one should only set high enough threshold. Note that by setting the threshold to infinity, errors are simply ignored and DW tables are always updated with newly arrived data.

A. The Fragmentation Schema

The first step in the process of horizontal fragmentation is to model the fragmentation scheme. Fragmentation scheme consists of the definition of a set of fragments that include all the attributes and tuples from the data warehouse. Horizontal fragmentation divides table into subsets of tuples. To be generic, the process of horizontal fragmentation must allow fragmentation per arbitrary filter. It should be implemented without changing the semantics of data warehouse and the resulting fragments must be disjoint.

According to [14] fragmentation is considered correct if it satisfies the following three rules:

(R1) Completeness: each tuple contained in the relation r , that is being decomposed into horizontal fragments r_1, r_2, \dots, r_n , can be found in one or more fragments r_i .

(R2) Reconstruction: for relation r that is being decomposed into horizontal fragments r_1, r_2, \dots, r_n , must be possible to define relational operator ∇ such that $r = \nabla r_i$, ($i=1, \dots, n$). The rule assures the reconstruction of relation without loss of information. The ∇ operator is usually the union operator: $r = r_1 \cup r_2 \cup \dots \cup r_n$.

(R3) Disjointness: horizontal fragments r_1, r_2, \dots, r_n of relation r are said to be disjoint if $r_i \cap r_j = \emptyset$, $\forall i, j \in \{1, \dots, n\}$, $i \neq j$.

Horizontal fragments are being defined using formula F .

Definition 1 Let r be the relation with schema R and let A and B be an attributes $A, B \in R$. Let θ be the relational operator from the set $\{=, \neq, <, \leq, >, \geq\}$ and let c be the constant from the domain of the attribute A . Formula F is defined by the following recursive expression:

1. $A \theta B$, $A \theta c$, $c \theta A$ are formulas. These formulas are called simple formulas.
2. If F_1 and F_2 are formulas, then $F_1 \wedge F_2$, $F_1 \vee F_2$, $\neg F_1$, $\neg F_2$ are formulas too.
3. Nothing else is a formula.

Definition 1 includes formulas applicable over domains with total order defined. In domains with no total order defined, only operators $=$ or \neq can be used for θ .

Definition 2: Let $\{F_1, F_2, \dots, F_n\}$ be the set of formulas defined on schema R of the relation r . Let for $\forall i, j \in \{1, \dots, n\}$ and each tuple $t \in r$ be satisfied $F_i(t) \wedge F_j(t) = \perp$ and $F_1 \vee F_2 \vee \dots \vee F_n(t) = T$.

Horizontal fragmentation transforms relation r with schema R into set of relations (fragments) f_1, f_2, \dots, f_n with schema R , such that $f_i = \sigma_{F_i}(r)$, ($i = 1, \dots, n$).

Such fragmentation meets each of the conditions: completeness, reconstruction and disjointness. Operator ∇ in this case is the union operator: $r = \cup f_i$, ($i = 1, \dots, n$).

With Definition 2, the so called primary horizontal fragmentation is defined. It is called primary because formulas that define horizontal fragments include only attributes from relations being fragmented. Derived horizontal fragmentation is carried out for relations associated with other fragmented (primary or derived) relations by means of foreign key.

B. ETL process with horizontal fragmentation

ETL process starts with the extraction phase of the source data (Figure 1), followed by the transformation and loading phase. In addition to metadata required for automated implementation of these phases of the ETL process, the data staging area contains a copy table for each source table as well as for each dimension and fact table of the data warehouse. Copies of the source tables are being loaded with source data in the extraction phase, while the remaining copies of the tables are being used in transformation and loading phases of the ETL process. They enable back-end testing and loading of data that, after integrity constraints have been checked, can be quickly copied to the dimension and fact tables in the data warehouse. These duplicate tables form a parallel set of tables with exactly the same schema as the schema of original tables in sources or in the data warehouse. They also have identical integrity constraints - for example if the fact table having scheme $F(a, b, c)$ references dimensional table d having schema $D(b, d, e)$ via attribute b , that means that in the data staging area exists tables f' and d' with schemas $F'(a, b, c)$ and $D'(b, d, e)$ such as F' references D' via an attribute b . In addition to the above tables, each data warehouse fact and dimension table in data staging area has its error-table used to store erroneous tuples. In addition to attributes from correspondent original table schema, the schemes of these error-tables include a timestamp attribute (the time when the tuple was forwarded to the error table) and a textual attribute comment (indicates the integrity rule violated for a particular tuple).

Figure 4 shows a sequence of activities taking place in the loading phase of the ETL process with horizontal fragmentation implemented. Besides a lifeline for the entire loading phase of the ETL process, separate lifelines are shown for the phases of ETL process related to dimension and fact tables in the data staging area (SA) and in the data warehouse (DW): "Dimensional Tables in SA", "Fact tables in SA", "Dimensional tables in DW" and "Fact tables in the DW". It can be seen from the

diagram that prior to transferring the data into the data warehouse all the necessary preparatory work is carried out in the data staging area. Activities taking place can be divided into several categories:

- dropping constraints
- loading data into tables
- determining horizontal fragments
- (re) creating constraints

Procedures with the DropConstraints or CreateConstraints prefix serve to drop or create integrity constraints (primary keys, foreign keys, unique and non-unique indexes). These procedures are performed in each

ETL process cycle, among other things, to speed up the loading phase in the data staging area as well as in the data warehouse (since it is faster to drop and recreate indexes than to adjust B⁺-trees on every single insert). Upon completion of the loading phase, integrity constraints and the indices are recreated. To facilitate the development and maintenance of ETL processes, these integrity constraints and indexes should be dropped and recreated programmatically. To be able to do so, it is necessary to store their definitions in the metadata repository.

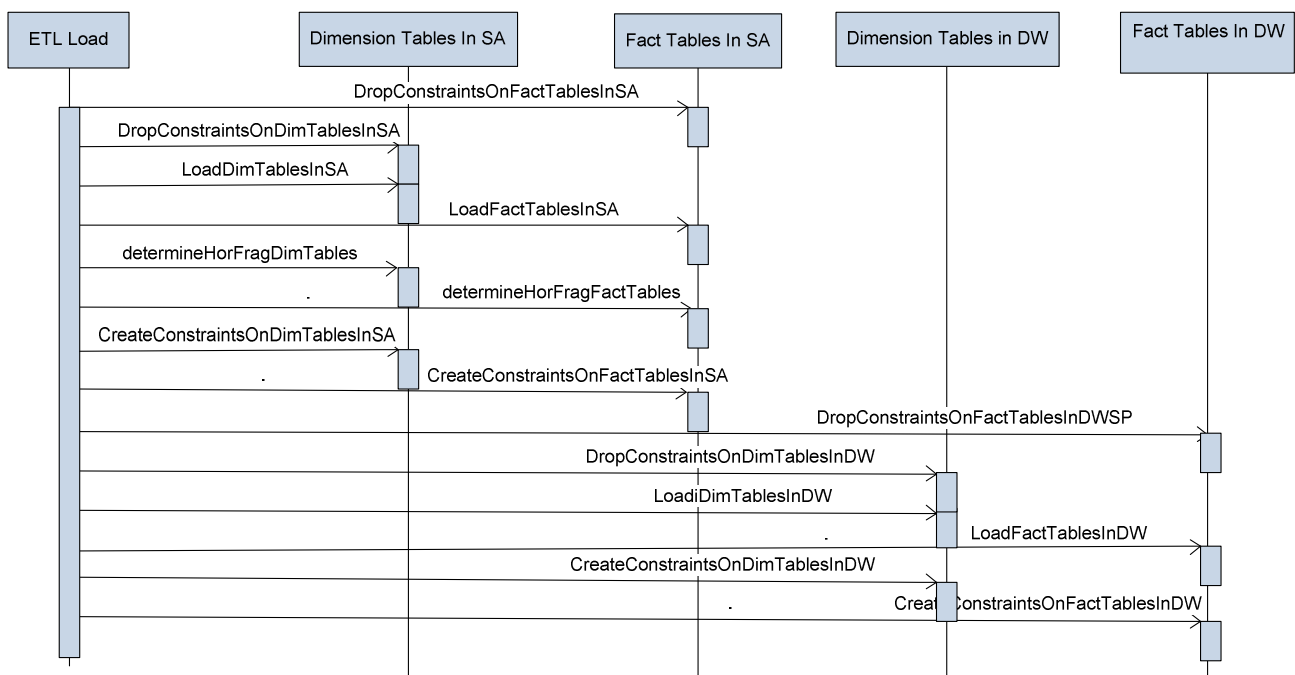


Figure 4 Activity sequence in loading phase of ETL process with horizontal fragmentation

After the constraints are dropped and data is loaded into the dimension and fact tables in the data staging area, we engage the procedure for finding tuples that do not meet the specific integrity constraint. Such tuples, if found, are forwarded (together with additional information about the error type and possible cause of errors) to corresponding error tables. It is possible that a certain source tuples in the extraction phase will not be copied from the sources to data staging area because they violate e.g. domain integrity or some other integrity rules. The absence of those erroneous tuples in data staging area may cause malfunction of other tuples that would otherwise be completely semantically correct. For example, due to violation of the domain integrity (e.g. BirthDate "1.1.1076" cannot be inserted into any SQL Server database), the tuple describing the student entity was not copied into the data staging area. On the other hand, tuples representing examinations of that student were successfully copied. However, missing student tuple causes reference integrity violation of the exam tuples.

Therefore those exam tuples cannot be loaded into the data warehouse. Finding and marking invalid tuples is an integral part of the cleansing and transformation phase of ETL process. Tuples marked as semantically incorrect will not be transferred into the data warehouse but forwarded to the appropriate error tables (and deleted from the originating tables). The process of finding and marking invalid tuples is carried out in data staging area only, prior to recreation of integrity constraints. After this procedure, each copy of dimensional and fact table in data staging area contains only valid tuples, while corresponding error tables hold the incorrect tuples.

Now, when all that is done, production fact and dimension tables can be loaded. During that process, the data warehouse will be unavailable to the users, so this must be done as quickly as possible. The proposed ETL process is designed with that in mind. The role of duplicate tables in the SA is primarily to reduce that time interval.

Again, loading is performed in three phases:

- constraints are dropped
- horizontal fragments are loaded
- constraints are recreated

The decision whether the horizontal fragment of the dimension or fact table in the data warehouse will be replaced with the corresponding horizontal fragment from the data staging area is based on a defined threshold. For example, suppose that there are two tuples in the exam table with the wrong examination date (e.g., set to the date ten years in the future) for the institution HEI89 from examples in Figure 2 and Figure 3. These two tuples will be detected by the ETL process and copied into the corresponding error table - errFExam. If the HEI89 institution defined the threshold for the fExam table to a value less than 2, that horizontal fragment will not be refreshed in that refresh cycle.

III. METADATA FOR THE HORIZONTAL FRAGMENTATION PROCESS

In order to be generic and easy to implement in the existing systems, horizontal fragmentation procedure is

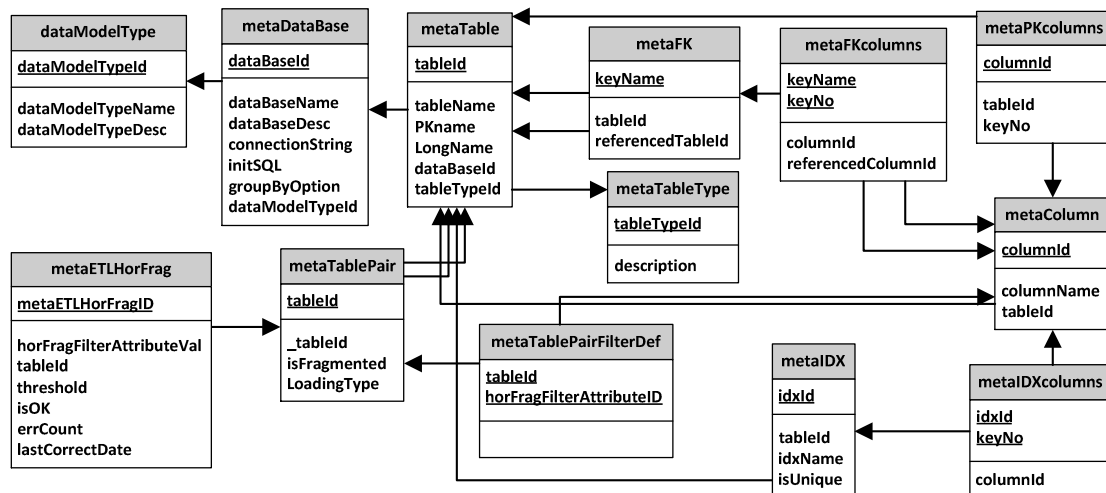


Figure 5 Metadata for the horizontal fragmentation process

The metaDataBase table is used for storing the data describing sub-systems of data warehousing system (source databases, staging area and data warehouse). For each fact and dimension table we introduce a matching “prime table” [1]. Prime tables form a parallel set of tables to the production tables having the same structure and analogous constraints; e.g. if a fact table $F(a, b, c)$ references dimension table $D(b, d, e)$ via attribute b , then $F'(a, b, c)$ references $D'(b, d, e)$ via attribute b . In short, prime tables are used to prepare and test data, and when ready, prime table data is quickly copied to the production tables. Data describing relations between each data warehouse table ($metaTablePair.tableId$) with its prime table ($metaTablePair._tableId$) is stored in the $metaTablePair$ table. For such pair, fragmentation attribute (or multiple attributes) is defined ($metaTablePairFilterDef$). $HorFragFilterAttribute$ is usually the id of the

designed and implemented with the help of the metadata. Metadata is stored in the extended data dictionary and used by generic ETL procedures. This way, to implement horizontal fragmentation, it is only required to populate metadata tables and to insert ETL procedure calls at the appropriate point in the existing ETL process.

Figure 5 shows the relational model of the metadata. Certain information must be entered manually (e.g., whether a table is a fact table or a dimension table - $metaTableType$ and $metaTable.tableTypeID$).

However, most of the information is already available in the DBMS data dictionary: tables, attributes, primary keys, foreign keys and indexes. This information can be copied automatically (different scripts have to be written for various DBMSs) to our metadata tables ($metaTable$, $metaColumn$, $metaPKColumns$, $metaFK$, $metaFKColumns$, $metaIDX$, $metaIDXColumns$). This information is used by the procedures (that generate and execute dynamic SQL) to drop and recreate integrity rules (foreign and primary keys) and indexes before and after data loading, and thus speed things up.

department, institution, etc. The $metaETLHorFrag$ table holds data about allowable number of incorrect tuples (threshold) for different fragments. We designed $horFragFilterAttributeVal$ attribute as a char field so that any data type values (or multiple values) could be stored in there. To reduce the number of tuples (thresholds) that user has to define, we adopt a convention that the threshold for the missing values (e.g. department) is default (in our case, zero). The $metaETLHorFrag$ table is, on procedure completion, updated with the outcome attributes ($isOK$, $errCount$, $lastCorrectDate$).

IV. IMPROVING THE COMPLETENESS AND TIMELINESS

In the context of horizontal fragmentation of the data warehouse tables, we do not consider the problems of data quality between the real world and the source data, but the data quality problems that arise from the extraction, transformation and loading of data from a

relational source into the data warehouse environment. Accordingly, the relational data source is considered complete and timely and incompleteness and timeliness appear in the data warehouse when it does not represent the exact image of the source data.

A. *Completeness*

The data warehouse is complete as much as the data which should be extracted from the data source, is, in fact, extracted, transformed and loaded into the data warehouse.

Let r be the source table comprising of M tuples with schema R comprising of N attributes. Let s be the data warehouse DW table having r as a source table. Let $ref(r)$ be the referent table for r consisting of all semantic correct tuples of the r , $card(ref(r)) \leq card(r)$. The completeness of the table s is defined with the following expression:

$$Q_{Compl}(s) = \frac{card(s)}{card(ref(r))}$$

Completeness of the table defined with above expression takes value from the interval $[0, 1]$. Completeness of the data warehouse DW comprising of K tables, according to [4] is defined by the completeness of the consisting tables:

$$Q_{Compl}(DW) = \frac{\sum_{k=1}^K Q_{Compl}(r_k) * g_k}{\sum_{k=1}^K g_k}$$

$g_k \in [0,1]$, in the above expression, represents the relative importance of the table r_k in the data warehouse DW , and provides the opportunity that the completeness of table, that is estimated more important, affects the completeness of the entire data warehouse more. The possible approach

[15] to determine the weight factor g_k of the r_k table takes into account the number of attributes of a relational schema (N_k) and the number of tuples in a table (M_k):

$$g_k = \frac{M_k * N_k}{\sum_{l=1}^K M_l * N_l}$$

An example from the Figure 6 and Figure 7 shows how the horizontal fragmentation process improves the completeness of the data warehouse table and consequently the completeness of the entire data warehouse. Consider an isolated segment of a data warehouse that integrates data from more higher education institutions. The observed segment consists of one dimension ($dStudent$) and one fact table ($fExam$). In the source tables ($student$ and $exam$) data from number of higher education institutions are stored as well. The numbers in cells are number of tuples for particular fragment and table. The data from three higher education institutions (HEI 1, HEI 2 and HEI 3) are taken into account in the example. Two consecutive cycles of data warehouse refreshment c_i and c_{i-1} are observed. In the c_i cycle there was no erroneous tuples in the source and the warehouse is updated with all the source data. Upon completion of the c_i cycle completeness of the data warehouse is maximal and equals 1. In the c_{i+1} cycle, for the higher education institution HEI 3, one erroneous tuple appears in the $student$ table and two erroneous tuples in the $fExam$ table. The ETL process without horizontal fragmentation implemented, due to erroneous tuples, will not refresh the content of $dStudent$ nor the content of $fExam$ table. The content of $dStudent$ and $fExam$ after the c_{i+1} cycle will be the same as after the c_i cycle completion.

	cycle c_i						cycle c_{i+1}					
	source			DW			source			DW		
	HEI 1	HEI 2	HEI 3	HEI 1	HEI 2	HEI 3	HEI 1	HEI 2	HEI 3	HEI 1	HEI 2	HEI 3
student/dStudent	4000	2000	1000	4000	2000	1000	4400	2200	1100	4000	2000	1000
exam/fExam	80000	40000	20000	80000	40000	20000	84000	42000	21000	80000	40000	20000

Figure 6 Completeness – data warehouse loading without horizontal fragmentation

1 erroneous tuple in student
2 erroneous tuples in exam

Completeness of the $dStudent$ and $fExam$ table is calculated with the following expressions:

$$Q_{Compl}(dStudent) = \frac{card(dStudent)}{card(ref(student))} = \frac{4000 + 2000 + 1000}{4400 + 2200 + (1100 - 1)} = \frac{7000}{7699} = 0,9092$$

$$Q_{Compl}(fExam) = \frac{card(fExam)}{card(ref(exam))} = \frac{80000 + 40000 + 20000}{84000 + 42000 + (21000 - 2)} = \frac{140000}{146998} = 0,9523$$

Completeness of the data warehouse segment is calculated based on the completeness of the tables contained in observed segment. Assuming that each table affects data warehouse completeness with an equal weight, weighting factor g_k equals to 1 for both tables.

$$Q_{Compl}(DW) = \frac{Q_{Compl}(dStudent) * g_1 + Q_{Compl}(fExam) * g_2}{1 + 1} = \frac{0,9092 * 1 + 0,9523 * 1}{2} = 0,9307$$

illustrates refreshment of the data warehouse segment using an ETL process with implemented horizontal fragmentation. A natural criterion for the fragmentation in this example is the identifier of the higher education institutions. Assuming that the tolerated number of erroneous records ($metaETLHorFrag.threshold$) is set to 0 for each institution and each table, the associated fragment of the $dStudent$ and $fExam$ for the HEI 3 in the c_{i+1} cycle will not be refreshed, while the fragments associated to the remaining higher education institutions having no errors will be refreshed.

The completeness of *dStudent* and *fExam* tables and the entire data warehouse, with horizontal fragmentation implemented, amounts to:

$$Q_{compl}(dStudent) = \frac{4400 + 2200 + 1000}{4400 + 2200 + (1100 - 1)} = \frac{7600}{7699} = 0,9871$$

$$Q_{compl}(fExam) = \frac{84000 + 42000 + 20000}{84000 + 42000 + (21000 - 2)} = \frac{146000}{146998} = 0,9932$$

$$Q_{compl}(DW) = \frac{0,9871 * 1 + 0,9923 * 1}{2} = 0,9901$$

	cycle ci						cycle ci+1					
	source			DW			source			DW		
	HEI 1	HEI 2	HEI 3	HEI 1	HEI 2	HEI 3	HEI 1	HEI 2	HEI 3	HEI 1	HEI 2	HEI 3
student/dStudent	4000	2000	1000	4000	2000	1000	4400	2200	1100	4400	2200	1000
exam/fExam	80000	40000	20000	80000	40000	20000	84000	42000	21000	84000	42000	20000

Figure 7 Completeness – data warehouse loading with horizontal fragmentation

1 erroneous tuple in student
2 erroneous tuples in exam

B. Timeliness

The data arrive into the data warehouse timely if it is updated (copied) within the first possible data warehouse refresh cycle.

Expression proposed in [15] [16] is used to define timeliness on the attribute value level:

$$Q_{Time}(v_o(t_i, A_j)) = \exp(-decline(A) * age(v_o(t_i, A_j)))$$

$Q_{Time}(v_o(t_i, A_j))$ denotes the probability that the attribute value is still valid. $v_o(t_i, A_j)$ is the observed and recorded value of the attribute A_j of the tuple t_i .

The *decline* (A) is the decline rate indicating how many values of the attribute considered become out of date on average within one period of time. The *age* ($v_o(t_i, A_j)$) denotes the age of the attribute value $v_o(t_i, A_j)$ which is computed by means of two factors: the instant when

attribute value timeliness is quantified and the instant of attribute value acquisition.

Assuming that the data warehouse is being refreshed once a day, decline and age of an attribute value are calculated with a day as the time unit. The example illustrated with Figure 8 and Figure 9 shows the improvement of the timeliness of attribute values, tuple, table, and the entire data warehouse accomplished with the horizontal fragmentation.

Analyzing changes in source tables for the HEIS [1] for regarded warehouse segment, we found that, on average, 2 % tuples change in one day in the student table and 0.6 % in the exam table. For a precise calculation of the timeliness dimension the decline of each relevant attribute should be assessed. However, that is not in the focus of this paper and so we consider here that the variability of all attributes is the same.

	cycle c _i						cycle c _{i+1}					
	source			DW			source			DW		
	HEI 1	HEI 2	HEI 3	HEI 1	HEI 2	HEI 3	HEI 1	HEI 2	HEI 3	HEI 1	HEI 2	HEI 3
student/dStudent	4000	2000	1000	4000	2000	1000	4000	2000	1000	4000	2000	1000
#of updated tuples							200	100	50	200	100	50
exam/fExam	80000	40000	20000	80000	40000	20000	80000	40000	20000	80000	40000	20000
#of updated tuples							400	200	100	400	200	100

Figure 8 Timeliness- data warehouse loading without horizontal fragmentation

1 erroneous tuple in student
2 erroneous tuples in exam

The timeliness of an arbitrary attribute of the student table with age of one and two days amounts to:

$$Q_{Time}(v_o(t_i, A_j)) = \exp(-0,002 * 1) = 0,9977$$

$$Q_{Time}(v_o(t_i, A_j)) = \exp(-0,002 * 2) = 0,9959$$

Similarly, the timeliness of an arbitrary attribute of the exam table amounts to:

$$Q_{Time}(v_o(t_i, A_j)) = \exp(-0,0006 * 1) = 0,9993$$

$$Q_{Time}(v_o(t_i, A_j)) = \exp(-0,0006 * 2) = 0,9987$$

If an attribute value fails to update (due to an error) at the very first cycle its age increases and timeliness decreases. The timeliness of the tuple is determined based on the timeliness of its attribute values. With the simplification that all attributes are equally important and have an equal decline rate and age, timeliness of tuple will be equal to the timeliness of the single attribute value. The timeliness of the r table comprising of the $card(r)$ tuples is determined by the timeliness of the its tuples:

$$Q_{Time}(r) = \frac{\sum_{j=1}^{card(r)} Q_{Time}(t_j)}{card(r)}$$

The timeliness of the data warehouse DW is determined by the timeliness of its tables:

$$Q_{Time}(DW) = \frac{\sum_{k=1}^K Q_{Time}(r_k) * g_k}{\sum_{k=1}^K g_k}$$

In the Figure 8 and Figure 9, the row titled “#of updated tuples” shows the number of tuples updated between two consecutive cycles c_i and c_{i+1} .

If the changes that had happened in the *dStudent* and *fExam* tables are successfully conducted in the c_{i+1} cycle, the timeliness of the tables and the data warehouse will be maximal. Suppose that the age of all tuples in the tables *dStudent* and *fExam* after the c_i cycle is one day. After the c_{i+1} cycle the age of altered (200 +100 +50 in *dStudent* and 400 +200 +100 in *fExam*) tuples is equal to one day and of all other tuples it is two days. However, in the presence of erroneous tuples, those changes will not reflect on *dStudent* and *fExam* in the c_{i+1} cycle if the strategy for the data

warehouse refreshment without horizontal fragmentation is employed.

The timeliness of *dStudent* and *fExam* after the c_{i+1} cycle is equal to:

$$Q_{Time}(dStudent) = \frac{\sum_{j=1}^{card(dStudent)} Q_{Time}(t_j)}{card(dStudent)} = \frac{(4000+2000+1000)*0.9959}{4000+2000+1000} = 0,9959$$

$$Q_{Time}(fExam) = \frac{(80000+40000+20000)*0.9987}{80000+40000+20000} = 0,9987$$

The data warehouse timeliness is equal to:

$$Q_{Time}(DW) = \frac{Q_{Time}(dStudent)*g_1 + Q_{Time}(fExam)*g_2}{1+1} = \frac{0,9959*1+0.9987*1}{1+1} = 0,9973$$

With implemented horizontal fragmentation changes (Figure 9) in the tuples belonging to the horizontal segments of HEI 1 and HEI 2 will be carried out successfully. After c_{i+1} cycle the age of modified tuples will be one, and of all other tuples two days.

The timeliness of *dStudent* and *fExam* and of the observed data warehouse segment amounts to:

$$Q_{Time}(dStudent) = \frac{(4000-200+2000-100+1000)*0.9959+(200+100)*0.9977}{4000+2000+1000} = 0,99597$$

$$Q_{Time}(fExam) = \frac{(80000-4000+40000-2000+20000)*0.9987+(4000+2000)*0.9993}{80000+40000+20000} = 0,99873$$

$$Q_{Time}(DW) = \frac{0,99597*1+0.99873*1}{1+1} = 0,99735$$

	cycle c_i						cycle c_{i+1}					
	source			DW			source			DW		
	HEI 1	HEI 2	HEI 3	HEI 1	HEI 2	HEI 3	HEI 1	HEI 2	HEI 3	HEI 1	HEI 2	HEI 3
student/dStudent	4000	2000	1000	4000	2000	1000	4000	2000	1000	4000	2000	1000
#of updated tuples							200	100	50	200	100	50
exam/fExam	80000	40000	20000	80000	40000	20000	80000	40000	20000	80000	40000	20000
#of updated tuples							4000	2000	1000	4000	2000	1000

Figure 9 Timeliness- data warehouse loading with horizontal fragmentation

1 erroneous tuple in student
2 erroneous tuples in exam

For determining timeliness exponential function (with negative exponent) with extremely slowly decreasing value is used. Hence, it is understandable that the differences in the timeliness of one day are not drastic – they are expressed in per mils. Never the less, the fact is that the process of horizontal fragmentation improves the timeliness of data in a given time interval.

In conclusion, the timeliness of the data warehouse refreshed employing ETL process without horizontal fragmentation will be the same as the timeliness of the same data warehouse refreshed employing ETL process with horizontal fragmentation when all the correct data are successfully updated. The benefit from horizontal fragmentation is in improving the timeliness in the occasions when, due to errors, correct data isn't updated.

V. CONCLUSION

In this paper we've discussed a generic procedure for horizontal fragmentation of the data warehouse. Proposed procedure is suitable when there are multiple somewhat independent departments (institutions, ...) that are generating data of the same structure (i.e. belonging to the same business process). For instance, different higher education institutions are generating students' exam records. In order to be generic, the procedure relies on metadata, which is presented and discussed. The criterion for the fragmentation of the table is arbitrary (e.g. institutionID). Using metadata, horizontal fragmentation can be customized to be more or less forgiving to errors. For instance, with high enough threshold values, algorithm recedes to a special case when all errors are ignored and data warehouse always gets updated with the new valid tuples (which is also a legit and popular approach).

We show that employing horizontal fragmentation improves data warehouse quality, namely, the completeness and timeliness data quality dimensions. Improvement is achieved in the cases when the ETL process without

horizontal fragmentation would not update correct data due to errors. The applicability of the proposed procedure has been tested in the real world project: data warehousing system integrating data from 90 higher educational institutions in the Republic of Croatia. Erroneous data detected in the process is logged and ultimately shown relevant users of the data warehousing system. Having ability to examine errors on-line, users are motivated to eliminate them and consequently raise the quality of data in the data warehouse, which presents a nice side benefit.

REFERENCES

- [1] R. Y. Wang, M. Reddy and H. B. Kon, "Toward quality data: An attribute-based approach," *Decision Support Systems*, pp. 349-372, 1995.
- [2] C. Batini and M. Scannapieca, *Data Quality Concepts, Methodologies and Techniques*, Heidelberg: Springer-Verlag Berlin, 2006.
- [3] L. Pipino, Y. W. Lee and R. Y. Wang, "Data Quality Assessment," *Communications of the ACM*, vol. 45, no. 4, 2002.
- [4] B. Heinrich, M. Kaiser and M. Klier, "Does the EU Insurance Mediation Directive help to improve Data Quality? - A metric-based analysis," in *Proc. of the 16th European Conference on Information Systems (ECIS)*, 2008.
- [5] M. Kaiser, "A Conceptual Approach to Unify Completeness, Consistency and Accuracy as Quality Dimensions of Data Values," in *European and Mediterranean Conference on Information Systems*, 2010.
- [6] L. Brkić, M. Baranović and I. Mekterović, "Improving the Completeness and Timeliness by Horizontal Fragmentation of Data Warehouse Tables," in *Proceedings of the 11th WSEAS International Conference on TELECOMMUNICATIONS and INFORMATICS (TELE-INFO '12)*, San Malo, 2012.
- [7] T. Jörg and S. Dessloch, "Formalizing ETL Jobs for Incremental Loading of Data Warehouses," in *Proceedings der 13. GI-Fachtagung für Datenbanksysteme in Business, Technologie und Web, Lecture Notes in Informatics*, Münster, Germany, 2009.
- [8] A. Behrend and T. Jörg, "Optimized Incremental ETL Jobs for Maintaining Data Warehouses," in *Proceedings of the Fourteenth International Database Engineering & Applications Symposium*,

2010.

- [9] P. Vassiliadis and A. Simitsts, "Near Real Time ETL," in *New Trends in Data Warehousing and Data Analysis, Annals of Information Systems*, R. W. S. Kozielski, Ed., Springer, 2008, pp. 19-49.
- [10] J. Guerra and D. A. Andrews, "Creating a Real Time Data Warehouse," Andrews Consulting Group, Inc., 2011.
- [11] C. Dell'Aquila, F. Di Tria, E. Lefons and F. Tangorra, "An Academic Data Warehouse," in *Proceedings of the 7th WSEAS International Conference on Applied Informatics and Communications*, Athens, Greece, 2009.
- [12] I. Mekterović, L. Brkić and M. Baranović, "Improving the ETL process of Higher Education Information System Data Warehouse," in *Proceedings of the 9th WSEAS International Conference on APPLIED INFORMATICS AND COMMUNICATIONS (AIC '09)*, 2009.
- [13] I. Mekterović, L. Brkić and M. Baranović, "Improving the ETL process and maintenance of Higher Education Information System Data Warehouse," *WSEAS transactions on computers*, vol. 8, no. 10, pp. 1681-1690, 2009.
- [14] M. Özsu and P. Valduriez, *Principles of distributed database systems*, 2nd ed., Upper Saddle River, New Jersey: Prentice-Hall International, Inc., 1999.
- [15] M. Kaiser, B. Klier and B. Heinrich, "How to Measure Data Quality?- A Metric-Based Approach," in *International Conference on Information Systems (ICIS)*, 2007.
- [16] B. Heinrich, "A Novel Data Quality Metric for Timeliness Considering Supplemental Data," in *17th European Conference on*

Information Systems (ECIS), Verona, 2009.

Ljiljana Brkić received her B.Sc, M.Sc and Ph.D. degree in Computer Science from the Faculty of Electrical Engineering and Computing University of Zagreb in 1992, 2004 and 2011 respectively. She has been affiliated with Faculty of Electrical Engineering and Computing as a research engineer from 1993 and as a research assistant at the Department of Applied Computing from 2006. Her research interests include data bases, data warehouses, business intelligence, information systems and programming paradigms.

Igor Mekterović received his B.Sc., M.Sc. and Ph.D. degree in Computer Science from the Faculty of Electrical Engineering and Computing, University of Zagreb in 1999, 2004 and 2008, respectively. Since 1999, he has been affiliated with Faculty of Electrical Engineering and Computing as a research assistant at the Department of Applied Computing. His research interests include databases, business intelligence and peer to peer systems.

Slaven Zakošek received his B.Sc., M.Sc. and Ph.D. degree in Computer Science from the University of Zagreb, Faculty of Electrical Engineering and Computing in 1987, 1992 and 2004 respectively. He has been affiliated with Faculty of Electrical Engineering and Computing as a research engineer, research assistant and assistant professor at the Department of Applied Computing. His research interests include data base systems and information systems.