

Fusionplex: resolution of data inconsistencies in the integration of heterogeneous information sources

Amihai Motro *, Philipp Anokhin

Department of Information and Software Engineering, George Mason University, University Drive, Fairfax, VA 22030-4444, USA

Received 4 December 2003; received in revised form 10 August 2004; accepted 9 October 2004

Available online 26 November 2004

Abstract

Fusionplex is a system for integrating multiple heterogeneous and autonomous information sources that uses *data fusion* to resolve factual inconsistencies among the individual sources. To accomplish this, the system relies on source *features*, which are meta-data on the merits of each information source; for example, the recentness of the data, its accuracy, its availability, or its cost. The fusion process is controlled with several parameters: (1) with a vector of feature weights, each user defines an individual notion of *data utility*; (2) with thresholds of acceptance, users ensure minimal performance of their data, excluding from the fusion process data that are too old, too costly, or lacking in authority, or numeric data that are too high, too low, or obvious outliers; and, ultimately, (3) in naming a particular fusion function to be used for each attribute (for example, *average*, *maximum*, or simply *any*) users implement their own interpretation of fusion. Several simple extensions to SQL are all that is needed to allow users to state these resolution parameters, thus ensuring that the system is easy to use. Altogether, Fusionplex provides its users with powerful and flexible, yet simple, control over the fusion process. In addition, Fusionplex supports other critical integration requirements, such as information source heterogeneity, dynamic evolution of the information environment, quick ad-hoc integration, and intermittent source availability. The methods described in this paper were implemented in a prototype system that provides complete Web-based *integration services* for remote clients.

© 2004 Elsevier B.V. All rights reserved.

Keywords: Information integration; Data inconsistency; Inconsistency resolution; Data fusion

1. Introduction and background

The *information integration problem* is defined as follows. Given a collection of heterogeneous and autonomous information sources, provide a system that allows its users to perceive the entire collection as a single source, query it transparently, and receive a single, unambiguous answer. *Heterogeneous* information sources are sources with possibly different data models, schemas, data representations, and interfaces.

Autonomous information sources are sources that were developed independently of each other, and are maintained by different organizations, that may wish to retain control over their sources.

An important issue in information integration is the possibility of *information conflicts* among the different information sources. These conflicts are at two different levels:

- *Intensional inconsistencies*: The sources are in different data models, or have different schemas within the same data model, or their data is represented in different natural languages or different measurement systems. Such conflicts have often been termed *semantic inconsistencies*.

* Corresponding author. Tel.: +1 7039931665; fax: +1 7039931638.
E-mail address: ami@gmu.edu (A. Motro).

- *Extensional inconsistencies*: There are factual discrepancies among the sources in data values that describe the same objects. Such conflicts are also referred to as *data inconsistencies*.

Extensional inconsistencies can only be observed after intensional inconsistencies have been resolved. That is, different attribute names in the schemas of different information sources must be mapped to each other and attribute values must be within the same measurement system, to conclude that these values indeed contradict each other.

Almost every system for information integration deals with intensional inconsistencies. In contrast, extensional inconsistencies have begun to receive their due attention only recently. This paper describes an information integration system that also resolves extensional inconsistencies.

1.1. The overall approach

Fusionplex is a development of an earlier information integration system, called Multiplex [19]. It is therefore a general data integration system, with support for information source heterogeneity, dynamic evolution of the information environment, quick ad-hoc integration, and intermittent source availability.¹ But the focus of Fusionplex is the resolution of extensional inconsistencies by means of *data fusion*.

The main principle behind Fusionplex is that “all data are not equal”. The data environment is not “egalitarian”, with each information source having the same qualifications (as assumed by Multiplex and most other systems). Rather, it is a diverse environment, in which information providers have their individual advantages and disadvantages. Some data are more recent, whereas other are more dated; some data come from authoritative sources, whereas other may have dubious pedigree; some data may be inexpensive to acquire, whereas other may be costlier. To resolve conflicts, Fusionplex looks at the qualifications of its individual information providers. Thus, it uses meta-data to resolve conflicts among data.²

Every Internet user is often confronted with the need to choose between alternatives: Which is the most trustworthy source? Which is the most reliable download site? Which is the least expensive newswire service? Some of these meta-data may be provided by the source itself (e.g., date of last update, cost), other meta-data may be obtained informally from other Internet users,

and there are also Web sites that are dedicated to calculating the quality of information and services provided by other sites (often through the evaluations of fellow users). So it is not far-fetched to assume that in the near future, given the Internet’s continuing, fast-paced expansion, such meta-data will become commonplace, possibly even in a standard format. In a more restricted information environment, comprising perhaps only a few dozen sources (possibly with a focus on a particular subject, such as business or medicine), it is quite conceivable that the multidatabase administrator will assign meta-data scores to its sources, and will keep updating these scores.

Our term for such meta-data is information *features*. Examples of features include: (1) *Timestamp*: The time when the information in the source was validated. (2) *Cost*: The time it would take to transmit the information over the network, or the money to be paid for the information, or both. (3) *Accuracy*: Probabilistic information that denotes the accuracy of the information. (4) *Availability*: The probability that at a random moment the information source is available. (5) *Clearance*: The security clearance level needed to access the information.

The other guiding principle of Fusionplex is that inconsistency resolution is a process in which users must be given a prominent voice. Depending on their individual preferences (which are subject to individual situations), users must be allowed to decide how inconsistencies should be resolved. Decisions are made in two phases: Some decisions are made ad-hoc at query-time, other decisions are more enduring and control all subsequent queries. One important query-specific decision is what constitutes “good” data. This decision is implemented by means of a vector of feature weights. In essence, this information constitutes this query’s definition of *data quality* and allows the system to rank the competing values according to their *utility* to the user. Other query-specific parameters are thresholds of acceptance with which users can ensure minimal performance of the data, excluding from the fusion process data that are too old, too costly, or lacking in authority. It also allows users to reject numeric data that are too high, too low, or obvious outliers. Yet another decision is the particular fusion function to be used for a particular attribute; for example, *average*, *maximum*, or simply *any*. Several simple extensions to SQL are all that is needed to allow users to state these resolution parameters. Altogether, Fusionplex provides its users with powerful and flexible, yet simple, control over the fusion process. For more details on this project, see [2]. Preliminary results were reported in [3].

The formal model of this work is the relational model. This model and the fundamental concepts of multidatabases are reviewed in Section 2. These fundamental concepts have been extended to include features. The extensions involve modifications to the basic relational

¹ These aspects are explained in detail in Section 2.3.

² A somewhat similar approach can be seen in the area of Internet search engines, where Google, one of the prominent search engines, weighs heavily the *authority* or *importance* of individual links in the ranking of its search results [13].

structures, the algebra, SQL, and the definition of multidatabases. They are discussed in Section 3.

To generate an inconsistency-free answer to a given query, Fusionplex must first retrieve all relevant data from the information sources. Section 4 describes how Fusionplex concludes which of the available information is relevant, and how this information is assembled in a “raw” answer. This intermediate product is termed *polyinstance* and it contains all the relevant data, including all possible inconsistencies. This polyinstance is the input to the resolution process. The first step in the resolution process is the identification of inconsistencies. In this process the tuples of the polyinstance are clustered in *polytuples*. The members of each polytuple are the different “versions” of the same information. Essentially, the inconsistency resolution process fuses the members of each polytuple in a single tuple. This fusion process is multiphased and is based on information provided either in the query itself or in the currently prevailing resolution policies. The detection and resolution of inconsistencies are the subject of Section 5.

The methods described in this paper were implemented in a prototype system called Fusionplex. The architecture and features of Fusionplex are described in Section 6, which also describes an experiment with a substantial example. Finally, Section 7 summarizes the contributions of this work and reviews several possible directions for further research. We begin with a brief review of research related to the subject of this paper.

1.2. Related research

Our discussion of related research is divided into three parts. We discuss first common approaches to the resolution of intensional (semantic) inconsistencies. As already mentioned, most systems for information integration address such inconsistencies. Because the area of information integration has been the subject of extensive research and development extending over a quarter of a century, this discussion does not attempt to be comprehensive. We then focus our attention on systems and approaches for dealing with extensional inconsistencies. We conclude this review with a discussion of information quality as it relates to the subject of this paper.

1.2.1. Intensional inconsistencies

Most systems for information integration create an overall description of the disparate information sources being integrated (the local sources). This global description is associated with the available descriptions of the local sources. Often, the global description is in the form of a database schema. Queries against the global schema are answered by fetching the appropriate information (as determined by the associations to the local sources). The global schemas are often conventional database

schemas, offering convenience to users, but different systems use different methods to store the global–local associations and to process global queries. Whatever their form, by their very nature, the global–local associations resolve semantic inconsistencies. Examples of this approach are the Information Manifold [17], SIMS [4], Infomaster [11], and Multiplex [19].

Instead of a global database schema, the global description can be a collection of *mediators*. Mediators are software modules that assemble global “objects” from the information provided by the local sources. Again, these software modules, created manually, resolve all semantic inconsistencies. Examples of this approach are TSIMMIS [10] and HERMES [34].

In our opinion, mediator-based systems suffer from an inherent weakness in that all integrated objects must be anticipated and pre-defined, and need to be redefined whenever the available information sources change. These problems are alleviated with the use of declarative languages for the definition of mediators, and software tools that assist in their construction, but they are not eliminated altogether. It can be argued that Multiplex (and other systems that do not use mediators, such as SIMS, the Information Manifold or Infomaster) offer more flexibility. A Multiplex query (a view of the global schema) may be considered a new “object” and its translation produces an ad-hoc “mediator”, describing how the global object is to be constructed from the presently available sources. The advantage of such “dynamic mediation” is that an unlimited number of global objects may be defined spontaneously. Moreover, whenever information sources change, only the global–local mapping needs to be updated: a single change as opposed to possibly numerous mediator changes.

1.2.2. Extensional inconsistencies

Whereas most systems for information integration resolve intensional inconsistencies, relatively few such systems address extensional inconsistencies: the discrepancies among the values obtained from different sources for the same data objects.

Systems that use mediators may attempt to resolve extensional inconsistencies in the mediators. For example, in HERMES, the mediator author must specify a conflict resolution method whenever a conflict might occur. In TSIMMIS, a mediator groups together information about the same real-world entity, possibly removing redundancies and resolving inconsistencies (the process is termed object fusion [29]). The weakness of mediators described earlier is also apparent in this form of inconsistency resolution. These predefined resolutions cannot be easily customized to individual users and specific queries.

Multiplex [19] both detects extensional inconsistencies and attempts to resolve them. However, inconsistency is approached at the “record level”: inconsistency occurs

when a global query results in two or more different *sets of records*. Multiplex then proceeds to construct an *approximation* of the true set of records, with a *lower bound* set of records (a *sound* answer) and an *upper bound* set of records (a *complete* answer). The two estimates are obtained from the conflicting answers through a process similar to voting. The lower bound set is contained in the upper bound set, and the true answer is estimated to be “sandwiched” between these two approximations. A significant limitation of this approach to inconsistency resolution is that Multiplex regards two records as describing entirely different objects, even if they are “almost identical” (e.g., identical in all but one “minor” field). Consequently, when two such records are suggested by two information sources, there is no attempt to recognize that these might be two descriptions of the same object (a task often referred to as object identification), and therefore no attempt to reconcile their conflicting values. The two records are simply both relegated to the upper bound estimate.

InFuse [8] is a comprehensive information integration system, whose capabilities are supported by the language FraQL [33]. FraQL is an extension of SQL with features for resolving an extensive variety of inconsistencies, both intensional and extensional. Extensional inconsistencies are resolved by reconciliation functions, which are special-purpose Java functions. While this solution can handle any type of reconciliation strategy, it requires a-priori programming, and possibly a different function for each and every conflict.

In a recent study [25], the possibility of resolving extensional inconsistencies using standard SQL is explored. The main advantage of this approach is that it uses existing capabilities for this new task. The approach has significant limitations, though, simply because SQL was never designed with this task in mind. For example, SQL’s aggregate functions are usually insufficient, and queries can become long and complex. Yet, this study may suggest to researchers the appropriate SQL extensions that are needed for this task.

Extensional inconsistencies have also been considered in the context of data cleaning. AJAX [9] is a tool for specifying data cleaning programs, either in a single source or in multiple sources. In particular, it incorporates features for clustering tuples that correspond to the same real-world object, and for merging the tuples of each such cluster in a single tuple; the latter is done by using a variety of aggregate functions.

Several approaches have attempted to resolve extensional conflicts based on the content of the conflicting data and possibly some probabilistic information that is assumed to be available. They either detect the existence of data inconsistencies and provide their users with some additional information on their nature (e.g., [1]), or they try to resolve such conflicts by returning a *probabilistic value*: a set of alternative values with attached

probabilities [7,35,18,5]. There is elegance in the probabilistic approach, because probabilistic values are more general than simple values, and the type of output of their resolution process is the same as the types of its input (probabilistic values). But the benefit of a probabilistic value to the database user is often in doubt. Another drawback is that probabilistic information must be provided for every data item in an information source. This is relatively rare, especially for Web-based sources. Fundamentally, rather than resolve inconsistencies by concluding *data* from the conflicting data values, probabilistic methods focus on concluding *probabilities* from the conflicting probability values. In other words, these approaches fuse probabilities not data, and can be viewed as managing *uncertainty*, rather than *inconsistency*.

A drawback common to most of these methods is that they do not provide their users with sufficient control over the inconsistency resolution process (which is often “hidden” in the definition of mediators). Users cannot influence the resolution process to arrive at the “best” answer, where “best” is defined by the particular user in a particular situation. Often, there are multiple ways to resolve extensional inconsistencies, and their suitability can only be judged by the user. For example, in one situation a conflict in a set of values may be deemed to be resolved best by choosing the most recent value; in another, with the same set of values, the most frequent value (the mode) may well be preferred.

Another common drawback of the methods discussed so far is that they disregard the fact that the information provided by their participating sources is often very different in its accuracy, reliability or availability. They ignore any such differences among their sources and simply assume that all sources are equally good.

1.2.3. Information quality

There has been increased awareness in the past decade of the importance of information quality for business information systems [38,36,37,12,28]. Much of the work has been on topics such as modeling data quality, improving data quality through quality control processes, and the effects of information quality on enterprise performance. The subjects relevant to this paper—quantitative assessment of the merits of information sources, and the use of such assessments in the process of information integration and inconsistency resolution—have received relatively little attention. We discuss here some of the more relevant works.

Two basic information quality criteria, soundness and completeness (collectively called information *goodness*), are the subject of [21,22,30]. The authors describe statistical methods for assessing these parameters, and how to use these assessments to (1) derive goodness estimates for answers to arbitrary queries, and (2) reconcile inconsistencies in multidatabase environments.

A comprehensive review of different types of quality meta-data and their assessment may be found [27] and [23, pp. 29–50]. The use of assessments to select preferred information sources in environments of multiple sources is the subject of [24].

An information integration method that uses information quality in the evaluation of global queries is described in [26,23]. The authors' departure point is a typical information integration system, which is based on a mediator architecture. In such a system, a query against a mediator results in multiple query plans against the underlying sources. The authors observe that information from different sources may be of different quality, resulting in query plans of different value to the user. They assume that each information source is annotated with various information quality parameters. This information is then used to derive the quality of the different plans generated in the process of query translation. Only the highest ranking query plans are used to materialize the answer to a query. In general, the use of quality meta-data in Fusionplex is similar to its use in [26,23], but there are some notable differences. Whereas the latter assumes that the meta-data are available during query translation and uses them for choosing the execution plan, Fusionplex retrieves the meta-data along with the data and uses them to clean up its answers according to a specification provided by the user. More importantly, the main purpose of [26,23] is to answer queries with data of high quality, whereas Fusionplex, after using the meta-data to select and rank the data it accumulates, applies a further process of identifying extensional inconsistencies and resolving them by means of fusion.

2. Multidatabase concepts

In this section we provide a brief overview of fundamental multidatabase concepts used in this paper. These concepts are adopted from [19].

2.1. Relational databases

The relational data model was adopted for this work. This choice was motivated by the fact that the relational model is widely used and standardized, most production-quality database management systems implement this model, and most of the information sources that require integration are relational. Our terminology is mostly standard [31]. Throughout the paper, we use both relational algebra and SQL notations to describe views and queries.

Our resolution methodology could introduce *null values* into relations. This requires appropriate extensions to the relational model to determine the results of comparisons that involve nulls. Codd's three-valued logic [6]

is adopted for this purpose. In this logic, comparisons that involve nulls evaluate to the value *maybe*. Different interpretations of such *maybe* values can be provided. In general, a *permissive* interpretation will map *maybe* values to *true*, and a *restrictive* interpretation will map *maybe* values to *false*. In each situation, the interpretation of choice will be stated.

2.2. Schema mappings

Consider a database (D, d) , where D is the database schema and d is its instance. Let D' be a database schema whose relation schemas are defined as views of the relation schemas of D . The database schema D' is said to be *derived* from the database schema D . Let d' be the database instance of D' which is the extension of the views D' in the database instance d . The database instance d' is said to be *derived* from the database instance d . Altogether, a database (D', d') is a *derivative* of a database (D, d) , if its schema D' is derived from the schema D , and its instance d' is derived accordingly from the instance d .

Let (D_1, d_1) and (D_2, d_2) be two derivatives of a database (D, d) . A view V_1 of D_1 and a view V_2 of D_2 are *equivalent*, if for every instance d of D the extension of V_1 in d_1 and the extension of V_2 in d_2 are identical. Intuitively, view equivalence allows one to substitute the answer to one query for an answer to another query, although these are different queries on different schemas.

Assume two database schemas D_1 and D_2 , that are both derivatives of a database schema D . A *schema mapping* (D_1, D_2) is a collection of view pairs $(V_{i,1}, V_{i,2})$, where $V_{i,1}$ is a view of D_1 , $V_{i,2}$ is a view of D_2 , and $V_{i,1}$ is equivalent to $V_{i,2}$, for every i .

As an example, the equivalence of attribute *Salary* of relation schema *Employee* in database schema D_1 and attribute *Sal* of relation schema *Emp* in database schema D_2 is indicated by the view pair

$$(\pi_{Salary}Employee, \pi_{Sal}Emp)$$

As another example, given the schemas *Employee* = $(Name, Title, Salary)$ in database schema D_1 , and *Manager* = $(Ename, Level, Sal)$ in database schema D_2 , the following view pair indicates that the retrieval of the salaries of managers is performed differently in each database:

$$(\pi_{Name, Salary} \sigma_{Title=manager} Employee, \pi_{Ename, Sal} Manager)$$

Schema mappings are instrumental in our definition of multidatabases.

2.3. Multidatabases

Assume that there exists a hypothetical database that represents the real world. This ideal database includes the usual components of schema and instance, which

are assumed to be perfectly correct. The relationships between actual databases and this ideal database are governed with two assumptions.

- *The schema consistency assumption (SCA)*. All database schemas are *derivatives* of the real world schema. That is, in each database schema, every relation schema is a view of the real world schema. The meaning of this assumption is that the different ways in which reality is modeled are all correct; i.e., there are no *modeling errors*, only *modeling differences*. To put it in yet a different way, all intensional inconsistencies among the independent database schemas are reconcilable.
- *The instance consistency assumption (ICA)*. All database instances are *derivatives* of the real world instance. That is, in each database instance, every relation instance is derived from the real world instance. The meaning of this assumption is that the information stored in databases is always correct; i.e., there are no *factual errors*, only different *representations* of the facts. In other words, all extensional inconsistencies among the independent database instances are reconcilable.

In this work we assume that the Schema Consistency Assumption *holds*, meaning that all differences among database schemas are reconcilable. These schemas are related through a *multidatabase* schema, which is yet another derivative of this perfect database schema. On the other hands, it is assumed that the Instance Consistency Assumption *does not hold*, allowing the possibility of irreconcilable differences among database instances. This means that the database instances are *not* assumed to be derivatives of the real world instance.

Formally, a *multidatabase* is

1. A *global* schema D .
2. A collection $(D_1, d_1), \dots, (D_n, d_n)$ of *local* databases.
3. A collection $(D, D_1), \dots, (D, D_n)$ of schema mappings.

The first item defines the schema of a multidatabase, and the second item defines the local databases in the multidatabase environment. The third item defines a mapping from the global schema to the schemas of the local databases. The schemas D and D_1, \dots, D_n are assumed to be derivatives of the real-world schema, but the instances d_1, \dots, d_n are not necessarily derivatives of the real-world instance (see Fig. 1). Note that there is no instance for the global database, and therefore a multidatabase is said to be a *virtual database*.

The “instance” of a multidatabase consists of a collection of global view extensions that are available from the local databases. Specifically, the views in the first position of the schema mappings specify the “contributed information” at the global level, and the views in the second position describe how these contributions are materialized.

As defined earlier, schema mappings allow to substitute certain views in one database with equivalent views in another database. In a multidatabase, the former database is the global database, and the latter is a local database.

The above definition of multidatabases provides four *degrees of freedom*, which reflect the realities of multidatabase environments.

First, the mapping from D to the local schemas is not necessarily *total*; i.e., not all views of D are expressible in one of the local databases (and even if they are expressible, there is no guarantee that they are mapped). This

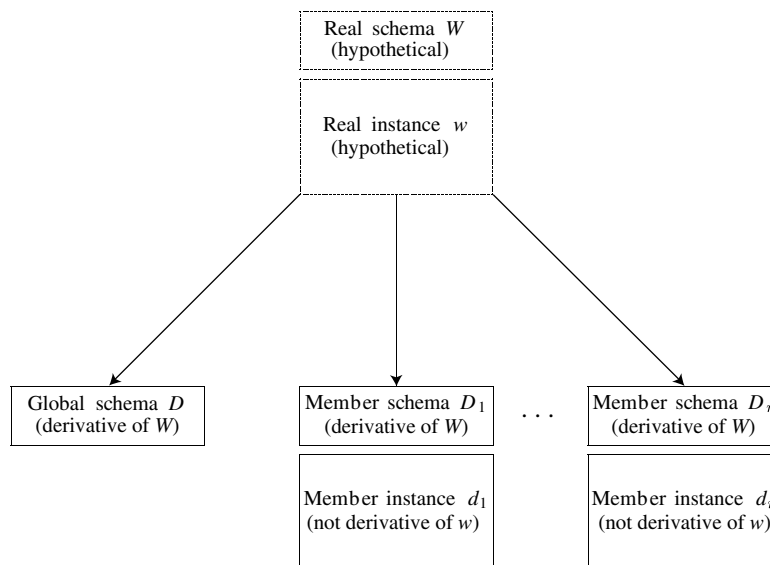


Fig. 1. Consistency assumptions in multidatabases.

models the dynamic situation of a multidatabase system, where some local databases might become temporarily unavailable. In such cases, the corresponding mappings are “suspended”, and some global queries might not be answerable in their entirety.

Second, the mapping is not necessarily *surjective*; i.e., the local databases may include views that are not expressible in D (and even if they are expressible, there is no guarantee that they are mapped). For example, a large database may share only one or two views with the multidatabase. This enables quick ad-hoc (i.e., not necessarily comprehensive) integration.

Third, the mapping is not necessarily *single-valued*; i.e., a view of D may be found in several local databases. This models the realistic situation, in which information is found in several overlapping databases, and provides a formal framework for dealing with multidatabase inconsistency. Since there is no assumption that the Instance Consistency Assumption holds, the local instances need not be derived from a single instance. Therefore, the fact that view pairs (V, V_1) and (V, V_2) participate in the schema mapping of a multidatabase does not imply that the extensions of V in the local databases are identical.

Fourth, while the definition appears to require that the local databases comply with the relational model, in practice they need not be relational, and the views in the second position of the schema mappings need not be relational model expressions. The only requirement is that they compute tabular answers. In other words, this definition allows for heterogeneity in the participating sources. The relational model provides only a convenient global description, and a communication protocol.

There have been many different models and architectures for information integration systems, and it is appropriate to place this architecture, which is at the basis of an earlier system Multiplex and the present system Fusionplex, in the proper context. In Section 1.1 we stated that most information integration systems offer a global description of the data environment and associate this global description with the descriptions of the local sources. Such systems can be classified by the type of their global–local associations. A classification offered in [14] distinguishes between architectures in which the local database schemas are defined in terms of the global schema (termed Local-as-View or LAV), and architectures in which the global schema is defined in terms of the local schemas (termed Global-as-View or GAV). Examples of the former type are SIMS, TSIMMIS and HERMES; examples of the latter type are the Information Manifold and Infomaster. The architecture of Multiplex is more powerful in that it associates views of the global schema with views of the local schemas. This hybrid approach earned the term GLAV.

3. Formal framework

To achieve our goal of resolving data conflicts, we introduce only one significant addition to the model described in Section 2, which we call *features*. This addition requires that we extend the notations for relational algebra queries (Section 3.2) and SQL queries (Section 3.3) as well.

3.1. Features

With the growth of the Internet, the number of alternative sources of information for most applications has increased enormously. To choose the most suitable source among the alternatives, users often evaluate information *about* the sources. These meta-data—whether provided by the sources themselves, by third-party sites dedicated to the ranking of information sources, or gained through prior experience—help users judge the suitability of each source for the intended use. Examples of such meta-data include:

- *Timestamp*: The time when the information in the source was validated.
- *Cost*: The time it would take to transmit the information over the network, or the money to be paid for the information, or both.
- *Accuracy*: Probabilistic information that denotes the accuracy of the information.
- *Availability*: The probability that at a random moment the information source is available.
- *Clearance*: The security clearance level needed to access the information.

These meta-data are referred to as source *features*. Each feature is associated with a domain of possible values, and a total order is assumed to be defined on the domain. For each information source that possesses a particular feature, a value from that domain is available. For example, the domain of *availability* could be the interval $[0, 1]$, the values of *cost* could range between 0 and M , where M is an arbitrary number, and the domain of *clearance* could be $\{top-secret, secret, confidential, unclassified\}$.

To facilitate comparisons between different feature values, all features are *normalized*. Each feature value is linearly mapped to a number in the interval $[0, 1]$. The mapping is done so that high feature values are always more desirable than low values; that is, the worst feature value is mapped to 0, and the best to 1. For example, higher *availability* value means higher probability of the source being available. Similarly, higher *timestamp* value means the data is more recent. But notice that higher *cost* value means the data is cheaper to obtain.

Every information source has a set of features associated with it, and in practice, different sources may have

different features. Therefore, a global set of features F is defined for the entire multidatabase, as the *union* of the features of the participating information sources. Each source feature set is then augmented to the features in F by adding null values for the features it does not possess. For example, assume the global set of features is $F = \{\text{timestamp}, \text{cost}, \text{availability}\}$. An information source with the current timestamp, zero cost and no data about availability would have the features $\text{timestamp} = 1$, $\text{cost} = 1$, $\text{availability} = 1$.

Our concept of features corresponds to that of *information quality criteria* in the information quality literature.³ An extensive list of such criteria may be found in [38]; a discussion of criteria in the context of information integration may be found in [26,23]. Our definition of features associates the same feature value with an entire information source. That is, in this work all features are assumed to be *inherited* by all individual tuples and all their attribute values. For example, in case of *timestamp*, it is assumed that the source-wide *timestamp* value is also the *timestamp* value for every attribute in every tuple of the source. Such features have been called *source-specific criteria* [26].

Clearly, this assumption is restricting as it implies that the data in every source are homogeneous with respect to every feature. However, in situations where an information source is heterogeneous with respect to its features, it should be partitioned into several disparate parts, homogeneous with respect to their features. These parts would consequently be treated as separate information sources. For example, a relation with two different *timestamp* values and two different *accuracy* values will be partitioned into four homogeneous parts. To assure that the number of partitions remains manageable, data with “close” feature values should belong to the same partition, and the partition should then carry feature values that represent the overall performance of its data; e.g., either an average value of some kind or a lower bound value. Partitioning sources into parts with similar feature values is described in [21,22,26].

We now update the definition of multidatabases offered in Section 2. The only change is in the definition of schema mappings. Recall that mappings were made of pairs, where the first element of each pair was a view of the global database D , and the other a view of a local database D_i . These pairs are now extended to triplets with the addition of the set of source features F . With this third element, every information source now provides its meta-data along with its data.

3.2. Extended relational model and algebra

To take advantage of feature meta-data in the processing of global queries, we extend the standard relational model. We offer extensions to the definitions of relation schemas, relation instances, and the relational algebra operations that manipulate these structures.

Each relation schema is extended with all the features in F . Assuming $F = \{F_1, \dots, F_k\}$, a relation schema $R = (A_1, \dots, A_m)$ is now extended to $R = (A_1, \dots, A_m; F_1, \dots, F_k)$. Correspondingly, the tuples in each relation instance are extended with the appropriate feature values. Recall that feature columns of source instances have the same value for all their tuples, and that *null* is an appropriate value for both database attributes and features.

The extension to the relational algebra consists of modifications to three basic operations: selection, projection, and Cartesian product. The other two basic operations, union and difference, remain unchanged.

The selection operation is extended to use an additional predicate. Besides the usual predicate ϕ for selecting tuples by their attribute values, it uses a predicate ψ for selecting tuples by their features. Let r be an (extended) relation instance of an (extended) relation scheme R . Then $\sigma_{\phi, \psi}(r)$ is the set of (extended) tuples from r that satisfy both predicates. Both predicates have two interpretations for handling null attribute values and null feature values: a *restrictive* interpretation in which comparisons to *null* evaluate to *false*, and a *permissive* interpretation in which they evaluate to *true*.

The extended Cartesian product *concatenates* the database values of the participating relations, but *fuses* their feature values. The fusion method depends on the particular feature. The *cost* incurred in creating each new tuple is usually the *sum* of the costs of the input tuples; the *availability* of the new tuple is the *product* of the availability values of the input tuples; the *timestamp* is the *minimum* of the input timestamps; and so on. Recall that higher feature values imply better performance. Thus, fusing two feature values with the minimum reflects a worst case approach intended to guarantee minimal performance of the combined information. Note that for the feature *cost*, where the sum is used, the resulting feature values have to be normalized to the range $[0, 1]$. This normalization is intricate, and for simplicity, our examples use the *average* cost instead.⁴ The fusion of feature values must also consider situations when one of the two feature values is *null*. Consider, for example, the feature *timestamp*. A null value may

³ We note that some authors use the term *quality criteria* (or quality parameters) to denote only intrinsic characteristics, such as accuracy or recentness, whereas others include also “external” characteristics, such as availability or cost. For this reason we prefer the more general term *information features*.

⁴ To justify, assume input costs c_1 and c_2 . A simple normalization is to divide their sum by the maximal cost possible. Since each input is normalized, that maximum is 2. Note that these new costs should be compared to each other, and not to the input costs.

Table 1
The evaluation of $\pi_{(Salary,Name);(0.3,0.7,0)}\sigma_{(EmpID=ID);(cost>0.6)}(R \times S)$

<i>r</i>						
<i>Salary</i>	<i>EmpID</i>		<i>time</i>		<i>cost</i>	<i>avail</i>
10,000	1002		1.0		0.5	<i>null</i>
50,000	1003		1.0		1.0	0.5
20,000	1001		0.8		0.2	<i>null</i>
<i>s</i>						
<i>ID</i>	<i>Name</i>		<i>time</i>		<i>cost</i>	<i>avail</i>
1002	<i>Johnson</i>		0.7		1.0	<i>null</i>
1002	<i>Johansen</i>		0.8		0.8	<i>null</i>
1001	<i>Nguyen</i>		1.0		1.0	0.1
1004	<i>Smith</i>		1.0		<i>null</i>	1.0
$t_1 = r \times s$						
<i>Salary</i>	<i>EmpID</i>	<i>ID</i>	<i>Name</i>	<i>time</i>	<i>cost</i>	<i>avail</i>
10,000	1002	1002	<i>Johnson</i>	0.7	0.75	<i>null</i>
50,000	1003	1002	<i>Johnson</i>	0.7	1.0	<i>null</i>
20,000	1001	1002	<i>Johnson</i>	0.7	0.6	<i>null</i>
10,000	1002	1002	<i>Johansen</i>	0.8	0.65	<i>null</i>
50,000	1003	1002	<i>Johansen</i>	0.8	0.9	<i>null</i>
20,000	1001	1002	<i>Johansen</i>	0.8	0.5	<i>null</i>
10,000	1002	1001	<i>Nguyen</i>	1.0	0.75	<i>null</i>
50,000	1003	1001	<i>Nguyen</i>	1.0	1.0	0.05
20,000	1001	1001	<i>Nguyen</i>	0.8	0.6	<i>null</i>
10,000	1002	1004	<i>Smith</i>	1.0	<i>null</i>	<i>null</i>
50,000	1003	1004	<i>Smith</i>	1.0	<i>null</i>	0.5
20,000	1001	1004	<i>Smith</i>	0.8	<i>null</i>	<i>null</i>
$t_2 = \sigma_{(EmpID=ID);(cost>0.6)}(t_1)$						
<i>Salary</i>	<i>EmpID</i>	<i>ID</i>	<i>Name</i>	<i>time</i>	<i>cost</i>	<i>avail</i>
10,000	1002	1002	<i>Johnson</i>	0.7	0.75	<i>null</i>
10,000	1002	1002	<i>Johansen</i>	0.8	0.65	<i>null</i>
$t_3 = \pi_{(Salary,Name);(0.3,0.7,0)}(t_2)$						
<i>Salary</i>	<i>Name</i>	<i>time</i>	<i>cost</i>	<i>avail</i>		
10,000	<i>Johansen</i>	0.7	0.7	<i>null</i>		

be any value in the interval $[0, 1]$. Assume now that the other value is α . The fused value (the minimum) is in the interval $[0, \alpha]$. Although more specific, this information is still represented by the value *null*. The fusion of feature values (including the necessary normalizations) is covered in detail in [20].

Finally, the projection operation is extended to perform an additional function. After the usual removal of attributes (note that feature columns cannot be removed), it also resolves inconsistencies in the resulting relation instance. The semantics of this additional function operation are discussed in Section 5. We note, however, that this function depends on a tuple of feature weights w (and on a set of prevailing resolution policies). Thus, $\pi_{S;w}(r)$ projects the instance r on the attributes S (and the features F), while purging it from all inconsistencies.⁵

The example in Table 1 illustrates these extensions to the relational algebra. The relations are $R = (Salary, EmpID)$ and $S = (ID, Name)$ and the features are *timestamp* (abbreviated *time*), *cost*, and *availability* (abbreviated *avail*). The example follows the construction of an answer to the query “Salaries and names of employees, where cost is not less than 0.5, and the importance of timestamp and cost are 0.3 and 0.7, correspondingly. The (extended) relational algebra expression for this query is $\pi_{(Salary,Name);(0.3,0.7,0)}\sigma_{(EmpID=ID);(cost>0.6)}(R \times S)$. In the interest of generality, the feature values in r and s are not uniform; this may be the case when these relations were created from different sources. Admittedly, details of the extended projection have not been provided yet, but note that it identifies *Johnson* with *Johansen*.⁶

⁵ The projection is a logical choice for handling the resolution of inconsistencies, as it already handles the removal of duplicates.

⁶ This step is revisited in Section 5.3.

3.3. Queries

For practical purposes, the relational algebra extensions defined earlier have also been added to SQL. The extended SQL query statement, shown below, introduces three new constructs.

```

select [restrictive] ...
from ...
where ...
using use_stmt
with weight_stmt
;

```

The **using** clause implements the new selection predicate ψ for specifying the desired features of the answer. In analogy with the **where** clause, which restricts the answer set with a condition on *attributes*, the **using** clause restricts the answer set with a condition on *features*. The syntax of the **using** clause is

```

use_stmt ::= feature_name comparison feature_value
           [and feature_name comparison feature_value]
           ...

```

The **with** clause provides for the specification of the feature weight tuple w used in the inconsistency resolution process. With this clause users assign weights to features. The features not mentioned in the **with** clause are assigned a weight of 0. The syntax of the **with** clause is

```

weight_stmt ::= feature_name as feature_weight
                [, feature_name as feature_weight]
                ...

```

Finally, when **restrictive** is present, all comparisons to null values, either in attribute columns (the **where** clause) or in feature columns (the **using** clause), evaluate to *false*; otherwise, they are *true*.

Expressed with the extended SQL statement, the previous query is

```

select Salary, Name
from R, S
where EmpID = ID
using cost > 0.6
with timestamp as 0.3, cost as 0.7
;

```

User queries are limited to relational algebra expressions that contain these relational algebra operations: (extended) projection, (extended) selection (without negation), (extended) Cartesian product, union, and difference.

4. Data collection and assembly

The most common approach to data integration is query translation (e.g., [17,4,19]). In this process, each user query, expressed in terms of the global schema, is translated to a query over the global views (the views that define the information available from the sources). The challenge in this translation is that some of the data requested in the user query may be in *none* of the available information sources (although it is described in the global schema), or it may be in *several* of them.

Most translation methods deal with the former problem by translating the given query to the maximal feasible query over the global views (i.e., a maximal subview of the original translation). The latter problem is addressed either by selecting one of the sources and ignoring the others (e.g., [4]), or by unifying the information from all the sources (e.g., [17]). In the following, we extend this query translation method to handle data inconsistencies.

Recall that each schema mapping is a set of triplets. We amalgamate the triplets of the different mappings in a single set, and denote each triplet (V, URL, f) , where V is a global view, URL is an expression that is used to materialize this view from one of the participating information sources, and f is the features of this source (i.e., a set of values for the features in the global set F). Each such triplet is called a *contribution* to the virtual database.

Like most translation methods, we restrict contribution views (the V of each triplet) to relational algebra expressions that include only projections, selections and joins (PSJ views). Additionally, we assume that they do not contain comparisons *across* the view relations.

Obviously, not all contributions are needed for every query. To determine the contributions that are relevant to a given query, the following two-step process is applied. First, the sets of attributes of the query and a contribution are intersected. If the intersection is empty, the contribution is deemed not relevant. Next, the selection predicates of the query and the contribution are conjoined. If the resulting predicate is not *false*, then the contribution is considered relevant to the query.

From each relevant contribution we derive a unit of information suitable for populating the answer to the query. Such units are termed *query fragments*. Intuitively, to obtain a query fragment from a contribution one needs to remove from it all tuples and attributes that are not requested in the query, and to add null values for the query attributes that are missing from the contribution.

Fig. 2 illustrates the construction of query fragments for a query Q from two relevant contributions: $C_1 = (V_1, URL_1, f_1)$ and $C_2 = (V_2, URL_2, f_2)$. The left part of the figure shows v_1 and v_2 , their intersection (the shaded area), and the “ideal” answer q (the dashed

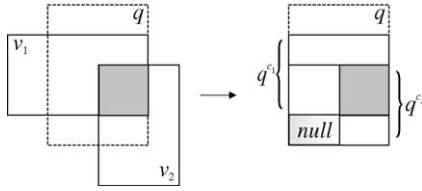


Fig. 2. Constructing fragments for query Q from contributions C_1 and C_2 .

box). The shaded rectangle represents an area of possible conflicts. The right part of the figure shows the two resultant query fragments, q^{C_1} and q^{C_2} . The area marked *null* contains null values as a result of V_2 not including all the attributes of Q .

At times, a source may provide values for one global attribute, yet through its definition, values of another global attribute may be inferred. Let V be a view. Assume that V includes the attribute A_i but excludes the attribute A_j , and assume that its selection predicate contains the equality $A_i = A_j$. Clearly, although this contribution does not provide the attribute A_j , it can be enhanced automatically to include it. Therefore, we add to V the attribute A_j , and we add to its instance v a column of A_j values which replicates the A_i values. This process, which is repeated for all equalities in V 's selection predicate, is called an *enhancement* of V and its result is denoted \bar{V} (the enhanced instance is denoted \bar{v}).

This discussion is summarized in this definition of query fragments:

1. Assume a query $Q = \pi_{(A_1, \dots, A_q)} \sigma_{\phi; \psi}(R_1 \times \dots \times R_p)$.
2. Assume a contribution $C = (V, URL, f)$. f is a tuple of feature values for the global feature set F (it has null values for features not available at URL).
3. Let \bar{V} be the enhanced contribution view, and let \bar{v} be its instance.
4. Denote $Y = \{A_1, \dots, A_q\} - \bar{V}$. Y is the difference between the schemas of the query and the enhanced contribution view. Let y be an instance of Y that consists of a single tuple t_{null} composed entirely of null values.
5. A *query fragment* derived from C for a query Q , denoted Q^C , is a view definition whose schema is $\{A_1, \dots, A_q\} \cup F$, and for every instance v of V , the instance of Q^C , denoted q^c , is $\sigma_{\phi; \psi}(\pi_{A_1, \dots, A_q}(\bar{v}) \times y \times f)$.

The concepts of contribution, enhancement and query fragment are illustrated in the following example. Assume a multidatabase relation with five attributes and three features: $R = (A, B, C, D, E; \text{timestamp}, \text{cost}, \text{availability})$, and consider a query

$$Q = \pi_{A, B, D, E} \sigma_{(C < 10); (\text{timestamp} > 0.5)} R$$

Further, assume a contribution $C = (V, URL, f)$ where

$$V = \pi_{A, B, C} \sigma_{B > 0 \wedge C = E} R$$

$$URL = \text{“http://www.aname.com/smth.cgi?action=retri-eve\&ID = 517”}$$

$$f = (\text{timestamp} = 0.7, \text{cost} = 0.8, \text{availability} = 1).$$

Finally, let the instance v of V be

A	B	C
1	2	7
2	2	11
3	4	4

First, this instance is extended to include the features

A	B	C	$time$	$cost$	$avail$
1	2	7	0.7	0.8	1.0
2	2	11	0.7	0.8	1.0
3	4	4	0.7	0.8	1.0

Next, because the selection predicate of V includes an equality $C = E$, V is enhanced to

A	B	C	E	$time$	$cost$	$avail$
1	2	7	7	0.7	0.8	1.0
2	2	11	11	0.7	0.8	1.0
3	4	4	4	0.7	0.8	1.0

Finally, the query fragment q^c formed from the contribution C is

A	B	D	E	$time$	$cost$	$avail$
1	2	<i>null</i>	7	0.7	0.8	1.0
3	4	<i>null</i>	4	0.7	0.8	1.0

Note that the column D contained in Q 's projection set is not available from the contribution, and is therefore presented as a column of null values.

From each relevant contribution a single query fragment is constructed. Some of these fragments may be empty. The union of all non-empty query fragments is termed a *polyinstance* of the query. Intuitively, a polyinstance encompasses all the information culled from the data sources in response to a user query.

From a user's perspective, a query against the virtual database is supposed to return a single consistent answer. By resolving all inconsistencies, the projection operation (to be described in Section 5) will convert this polyinstance to a regular instance.

Recall that query translation is a process in which a query over the global relations is translated to a query over the global views. Since each query fragment is a view over a global view, the polyinstance is a view over the global views as well. Hence, the construction of the polyinstance is simply a query translation process. Note that in the absence of data conflicts (i.e., when the ICA holds), this polyinstance is equivalent to the output of a conventional query translation algorithm.

5. Data inconsistency detection and resolution

An extensional inconsistency exists when two objects (or tuples in the relational model) obtained from different information sources are identified as *versions* of each other (i.e., they represent the same real-world object) but some of the values of their corresponding attributes differ. Note that such identification is only possible when both schema incompatibilities and representation differences, collectively referred to here as intensional (semantic) inconsistencies, have been resolved. The process of data integration requires two steps: inconsistency *detection* and inconsistency *resolution*.

5.1. Data inconsistency detection

Data inconsistency detection begins by identifying tuples of the polyinstance that are versions of each other. Several techniques have been suggested for identifying tuples or records originating from multiple sources [15,32,16]. The work described here is applicable to any of these methods. For simplicity, we assume that identification is by *keys*.

We assume that each global relation is fitted with a key. Subsequently, to find the tuples in the polyinstance that are versions of each other, one needs to construct the key of the answer and use it to cluster the polyinstance. The resulting clusters are termed *polytuples*. Each polytuple may be visualized as a table:

<i>ID</i>	<i>Name</i>	<i>Age</i>	<i>Salary</i>	<i>time</i>	<i>cost</i>	<i>avail</i>
5218	<i>Smithson</i>	38	75,000	0.8	0.5	1.0
5218	<i>Smith</i>	35	<i>null</i>	0.7	0.2	<i>null</i>
5218	<i>Schmidt</i>	35	77,000	0.7	0.8	1.0

It is possible to improve the efficiency and accuracy of polytuple clustering, by considering only “horizontal slices” of the polyinstance in which there is possible contention (recall the shaded areas in Fig. 2). These slices may be determined from the selection predicates of the individual contributions. For example, if ϕ_1 and ϕ_2 are the selection predicates of two contributions, then only the slice defined by $\phi_1 \wedge \phi_2$ might have conflicts. In other words, the membership of polytuples cannot span across different slices. For brevity, we do not address this refinement here, and additional details may be found in [2].

Consider now a contribution $C = (V, URL, f)$ and the view instance v materialized from this contribution. There is always a possibility that v is not an acceptable instance of V . For example, V may be a Cartesian product of two global relations, but the set of tuples materialized from the URL could not possibly be a Cartesian product. As another example, V may involve a selection $A = a$, yet v includes in its column A values different than a . In this paper, a contribution v that contradicts its definition V is ignored.

Once extensional inconsistencies have been detected, they need to be resolved; i.e., every polytuple should be reduced to a single tuple. This process is performed in two passes. In the first pass (Section 5.2), the members of each polytuple are ranked and purged according to user-specified preferences. In the second pass (Section 5.3), in each polytuple, in each attribute, remaining values are purged and then fused to a single value. Similarly, in each polytuple, in each feature, different values are fused to a single value. This procedure provides semantics to the extended projection operation, introduced in Section 3.2.

5.2. Utility function

Recall that the extended projection requires the specification of a tuple of feature weights w . With these weights, users prescribe the relative importance of the features in the resolution process. To use this information, a *utility function* is calculated for each member of every polytuple. Assume the user assigns weights w_1, \dots, w_k to the features F_1, \dots, F_k . Then a member with feature values f_1, \dots, f_k receives utility $u = \sum_{i=1}^k w_i f_i$. These utility values are used to *rank* the members of each polytuple. Using a pre-defined *utility threshold*, members of insufficient utility are discarded. Utility is calculated using only the features that are non-null for all members of the polytuple.

Since a polytuple may have several members of acceptable utility, this process is not sufficient for resolving polytuples. Actual resolution is achieved in a second pass, described next.

5.3. Resolution policies

The resolution of inconsistencies among different values can be based either on their features, such as *timestamp*, *cost* or *availability* (*feature-based* resolution policies), or on the data themselves (*content-based* policies). Ideally, a conflict resolution policy should be provided whenever data inconsistency is possible. However, to keep the number of policies under control, a policy is defined for each global attribute.

Within each polytuple, data inconsistency is resolved in each attribute separately (i.e., each column of values in the polytuple is fused to a single value). This process consists of two steps. First, some of the values in the given attribute are *eliminated*, based on attribute or feature values. Then, the remaining values are *fused* by a function, resulting in a single value. Correspondingly, a resolution policy is defined as a sequence of *elimination functions*, followed by a *fusion function*.

Elimination functions are either content-based or feature-based. Examples of elimination functions are *min* and *max*. In the previous example of a polytuple, consider the attribute *Salary*. Possible eliminations include

$max(timestamp)$, $max(availability)$, $min(cost)$, and $max(Salary)$. The former three are feature-based, whereas the latter is content-based. But other functions may also be used; for example, *above_average*, *top_five_percent* and *within_standard_deviation_of_the_mean*. Each function is applied in its turn (according to its place in the sequence) to the corresponding column. However, this step can still result in multiple values (e.g., several *Salary* values may share the maximal availability). Such multiple values are handled subsequently by the fusion function.

Fusion functions are always content-based. The fusion function is applied to the values of the attribute and to the values of the features, resulting in a single resolved value for each of them. Examples of fusion functions are *any* and *avg*. Consider the attribute *Salary* from the previous example. Applying *any* gives quick resolution by choosing a value at random, for example 77,000. Applying *avg* results in the *Salary* value 76,000. Functions other than *any* or *avg* may also be used. For example, *mode*, *average_without_extreme_values*, or any other function of the conflicting values. A valuable discussion of fusion functions can be found in [25].

When every attribute in the polytuple has been resolved, a simple tuple is obtained. Now the feature values of this tuple must be determined. Each feature value must reflect the feature values of all the participating polytuple members. A polytuple member is a participant in the final tuple, if it contributed either (1) a value to the final tuple, or (2) an input value to a fusion function (e.g., *avg*). The fusion of features is done as in the (extended) Cartesian product (Section 3.2).

Consider the final step in the example of Table 1. The utilities of the two tuples in t_2 are $0.3 \cdot 0.7 + 0.7 \cdot 0.75 = 0.735$ and $0.3 \cdot 0.8 + 0.7 \cdot 0.65 = 0.695$, and assume that both are above the threshold. Suppose that for *Name* the elimination function is $max(timestamp)$ and the fusion function is *any*, and suppose that for *Salary* an elimination function is not specified and the fusion function is *avg*. The final tuple is therefore (10,000, Johansen). Both input tuples contributed to this final tuple, hence its feature values are (0.7, 0.7, null).⁷

Altogether, the reduction of each polytuple to a simple tuple, reduces the polyinstance to a simple instance. This set of tuples is then presented to the user as an inconsistency-free answer to the query.

To give users full control over the process of inconsistency resolution, we provide a powerful and flexible resolution statement. This statement implements the full set of features discussed in this section:

```

for                 $A_i$ 
[keep[restrictive]  $e_1(F_1), \dots, e_n(F_n)$ 
fuse               $f$ 

```

Here, A_i is a global attribute name, e_1, \dots, e_n is a sequence of elimination functions for this attribute, and f is a content-based fusion function. Each F_i is a feature, and $e_i(F_i)$ indicates feature-based elimination (e.g., $max(timestamp)$ eliminates all but the most recent value). If F_i is not given, the elimination is content-based (e.g., $min()$ retains the smallest value). Elimination functions are applied in the order in which they appear in the **keep** clause. Note that the entire **keep** clause is optional. Frequently, multiple attributes would share the same resolution policy. To facilitate the specification, the resolution statement allows multiple attributes in its **for** clause: **for** A_{i_1}, \dots, A_{i_m} .

The handling of null values during the elimination phase is controlled by the **restrictive** keyword. The decision whether the value *null* satisfies an elimination function (either content-based or feature-based) is similar to the decision on *null* comparison in the **using** and **where** clauses: If the keyword **restrictive** is present, a *null* is assumed to not satisfy the elimination function; otherwise, it is assumed to satisfy it. Regardless, in the fusion phase, attribute values that are *null* are discarded, under the common assumption that, whenever available, non-null values should be preferred (if all the attribute values are *null*, then the fusion value is *null*). When determining the features of the fusion value, only the features of non-null attribute values are considered. The fusion of features is performed as explained in the definition of the (extended) Cartesian product. If any of the feature values is *null*, then their fusion is *null*.

A resolution statement is required for each attribute of the global schema. These statements may be supplied by the system, by domain experts, or defined by users at run-time. If an attribute is left without a resolution statement, a default policy is applied.

The following three examples illustrate different kinds of resolution policies. The examples all use the query

$\pi_{Name, Age, Salary} \sigma_{Position=Manager} Employee$.

and the polytuple

ID	Name	Age	Salary	time	cost	avail
5218	Smithson	38	75,000	0.8	0.5	1.0
5218	Smith	35	null	0.7	0.2	null
5218	Schmidt	35	77,000	0.7	0.8	1.0

1. *Content-based policy*. Consider this resolution policy that chooses any *Name*, the average *Age* and the minimal *Salary*:

⁷ The three feature values were obtained by using minimum, normalized sum (average), and product, respectively.

for	<i>Name</i>			fuse	<i>any</i>
for	<i>Age</i>			fuse	<i>avg</i>
for	<i>Salary</i>	keep	<i>min()</i>	fuse	<i>any</i>

This policy uses only attribute values for elimination. Its result is

<i>ID</i>	<i>Name</i>	<i>Age</i>	<i>Salary</i>	<i>time</i>	<i>cost</i>	<i>avail</i>
5218	Smith	36	75,000	0.7	0.5	<i>null</i>

- Feature-based policy.* Consider this resolution policy that chooses the *Name* that is most recent, and the *Age* and *Salary* that are least costly:

for	<i>Name</i>	keep	<i>max(timestamp)</i>	fuse	<i>any</i>
for	<i>Age</i>	keep	<i>max(cost)</i>	fuse	<i>avg</i>
for	<i>Salary</i>	keep	<i>max(cost)</i>	fuse	<i>any</i>

Note that least costly translates to maximal cost. This policy uses only feature values for elimination. Its result is

<i>ID</i>	<i>Name</i>	<i>Age</i>	<i>Salary</i>	<i>time</i>	<i>cost</i>	<i>avail</i>
5218	Smithson	35	77,000	0.7	0.65	1.0

- Mixed policy.* Consider this resolution policy that chooses the *Name* that is least recent, the lowest *Age* and the *Salary* that is least costly:

for	<i>Name</i>	keep	<i>min(timestamp)</i>	fuse	<i>any</i>
for	<i>Age</i>	keep	<i>min()</i>	fuse	<i>avg</i>
for	<i>Salary</i>	keep	<i>max(cost)</i>	fuse	<i>any</i>

This policy uses both attribute values and feature values for elimination. Its result is

<i>ID</i>	<i>Name</i>	<i>Age</i>	<i>Salary</i>	<i>timestamp</i>	<i>cost</i>	<i>avail</i>
5218	Schmidt	35	75,000	0.7	0.5	<i>null</i>

5.4. The overall resolution procedure

Our inconsistency resolution methodology is summarized in the following procedure. The procedure begins with a polyinstance: the “raw” answer, comprising the union of the query fragments that are extracted from the information sources, clustered into polytuples. Its output is a simple instance: the final, inconsistency-free answer to the query.

Input:

- A feature-based selection predicate ψ , obtained from the query’s **using** clause.
- A tuple of feature weights w , obtained from the query’s **with** clause.
- A utility threshold $0 \leq \alpha \leq 1$ to be used in conjunction with the weights.
- A resolution policy for every attribute of the global schema.

The first two items are specified by the user as part of the query. The last two items are assumed to be pre-defined (e.g., provided by either administrators or users).

Procedure:

- In each polytuple, remove members that do not satisfy the predicate ψ .
- In each polytuple, calculate the utility of each remaining member: $u = \sum_{i=1}^k w_i f_i$, and rank the members by their utility. Let u_0 denote the highest utility in the polytuple. Discard members whose utility is less than $\alpha \cdot u_0$.
- In each attribute of the polytuple, apply the resolution policy for that attribute: Eliminate attribute values according to the **keep** clause, and fuse the remaining values according to the **fuse** clause.
- Calculate the feature values of each resulting tuple, taking into the account the feature values of the polytuple members that contributed to this tuple.
- The tuples thus obtained for each polytuple comprise the final answer to the query.

Unless the keyword **restrictive** is specified in the query, polytuples may have members with null feature values. In such cases, features that include nulls are ignored, so that ranking is based on the features that are available for *all* members. For example, assume the utility function $u = 0.3 \cdot \text{timestamp} + 0.5 \cdot \text{cost} + 0.2 \cdot \text{availability}$ and the previous polytuple. The non-null features are *timestamp* and *cost*, and the utility function is changed to $u = 0.3 \cdot \text{timestamp} + 0.5 \cdot \text{cost}$. If there are null values in each of the features of a particular polytuple, then its members cannot be ranked reliably, and Step 2 is skipped for that polytuple.

If a global attribute is non-numeric (i.e., of type *string*), then the choice of possible elimination and fusion function is more limited. Some of the common functions may have to be specifically defined to operate on non-numerical values (for example, *min* or *max* may use lexicographical order). Other functions become meaningless; for example, *avg* is unlikely to be used as a fusion function for non-numeric attributes.

Finally, there could be situations in which users may wish to leave inconsistencies unresolved or resolved only in part. In general, three solutions are possible:

- No resolution.* When none of the resolution steps are applied, the answer is a polyinstance (a set of polytuples). From a theoretical point of view, such answers could violate key constraints (analogous to retention of duplicates).
- Pruning of polytuples.* This involves only the application of the **using** and **with** clauses, which remove tuples that either do not satisfy the feature selection

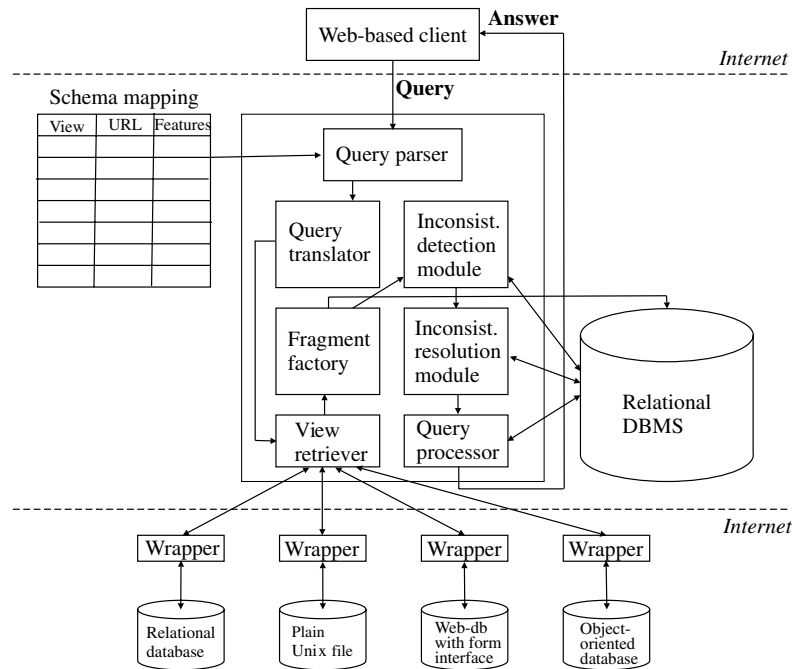


Fig. 3. Architecture of Fusionplex.

predicate, or are below the utility threshold (i.e., Steps 1 and 2 of the procedure). Such answers could still violate key constraints.

3. *Selective attribute resolution*. This involves leaving some (or all) attributes unresolved. It can be implemented with a special *do-not-resolve* policy. From a theoretical point of view, such answers could violate the first normal form, which dictates that each tuple has single value for each of its attributes.⁸

6. Implementation and experimentation

6.1. Implementation

The solutions presented in this paper were implemented in a prototype system. The system, called Fusionplex, is described in this section. Fig. 3 illustrates the overall architecture of Fusionplex.

Fusionplex conforms to a server–client architecture. The server is implemented in Java and contains the core functionalities described in this paper. At startup time, the server reads all configuration files, caches all source descriptions and creates temporary tables in a relational Database Management System. Then it starts listening for incoming connections from clients. Clients connect

to the server using simple line-based protocol. Each client passes to the server the name of a database that the client wishes to query and the query itself. The server processes the query and returns its result to the client, which formats it and delivers it to its user.

The core of Fusionplex consists of seven functional blocks:

1. The **query parser** parses the user query, checks its syntax and ensures that the relation and attribute names mentioned in the query are valid in the current virtual database.
2. The **query translator** determines the source contributions that are relevant to the given query by testing the intersection of the selection conditions of the query and of each contribution. It then calls the view retriever with the specifications of the contributions that were found relevant.
3. The **view retriever** consults the schema mapping and retrieves the relevant view instances from their corresponding URLs. It then attempts to enhance each instance with any attributes that participate in the view selection condition as part of an equality.
4. The **fragment factory** constructs the query fragments from the enhanced views and stores them in the relational database management system.
5. The **conflict detection module** assembles a polyinstance of the answer from the fragments, determines the possible areas of inconsistency by examining the selection predicates associated with the fragments and constructs the polytuples.

⁸ Note that using the *do-not-resolve* policy for every attribute is different from Solution 2.

6. The **conflict resolution module** resolves data conflicts in each polytuple according to the appropriate resolution policies. First, each mono-attribute polytuple is resolved with a single value. Then the features of the resolved tuples are determined.
7. The **query processor** processes the union of all the resolved tuples, by applying any aggregation and ordering specified in the query, and returns the query result.

Fusionplex also provides a client with a graphic user interface (GUI). This web-enabled client supports both direct input of queries as simple text, and guided query construction through the use of its Query Assistant. The Query Assistant allows users to create queries using an intuitive visual interface. The client consists of a set of CGI-scripts written in Perl and allows for minimal correctness check at the client side. Fig. 4 shows the client interface with an ongoing Query Assistant session.

As soon as the interaction with the Query Assistant is completed, a query in the SQL-like query language of Fusionplex is constructed and displayed in the query window. When the *Submit* button is pressed, this query is transmitted to the server, and the results returned from the server are displayed by the client (Fig. 5).

Fusionplex also incorporates a database management tool for defining and maintaining virtual databases. Authorized users can create new virtual databases, and modify existing ones. To create a new virtual database, the user must define its global relations and plug-in any number of contributions from existing information sources. In each contribution, the user must specify a global view, a matching URL, and the appropriate source features. Existing virtual databases can be modified by adding or removing relations or contributions.

To experiment with the Fusionplex system, several information sources were constructed. For the purpose

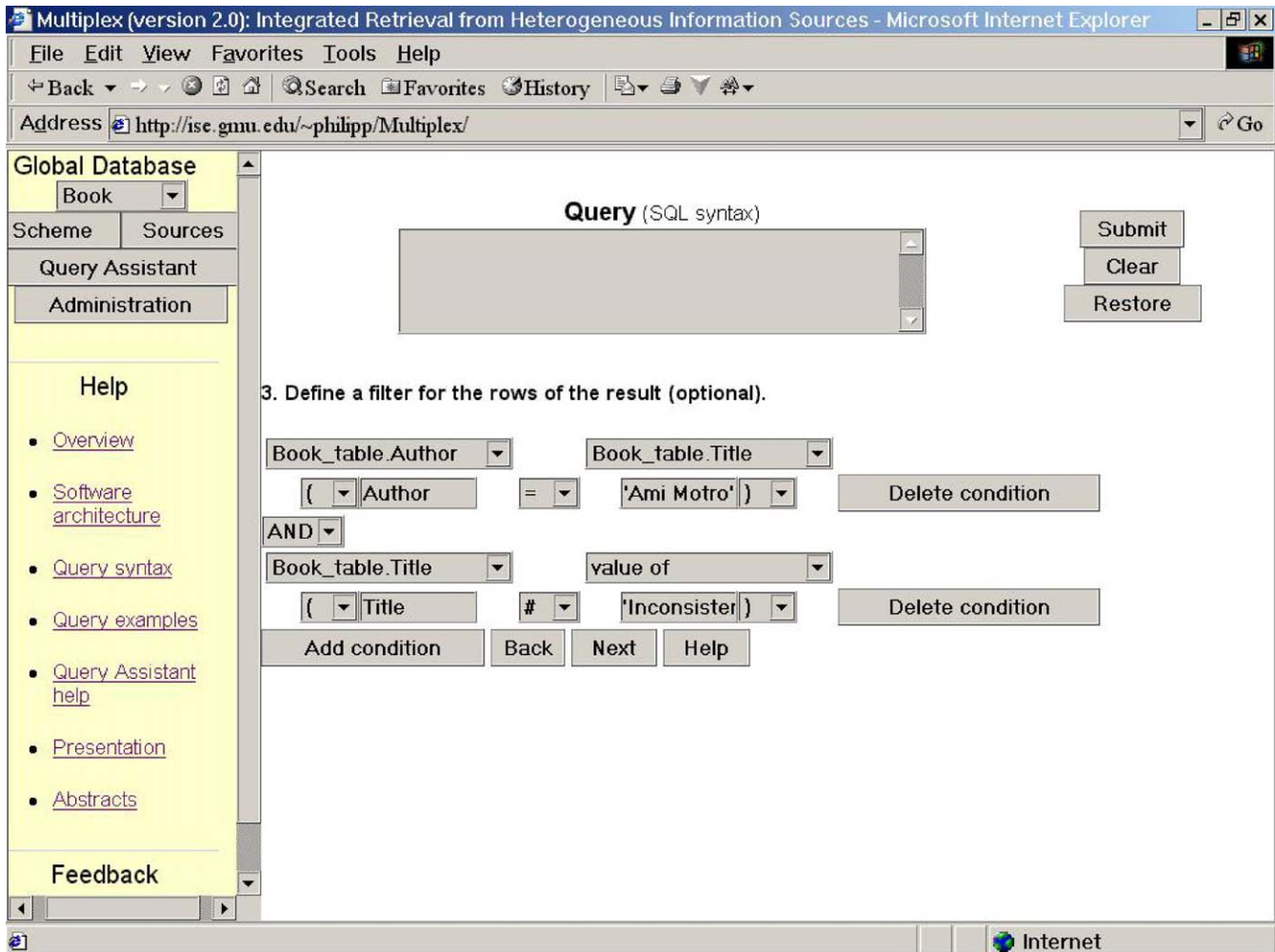


Fig. 4. The Fusionplex GUI and a Query Assistant session.

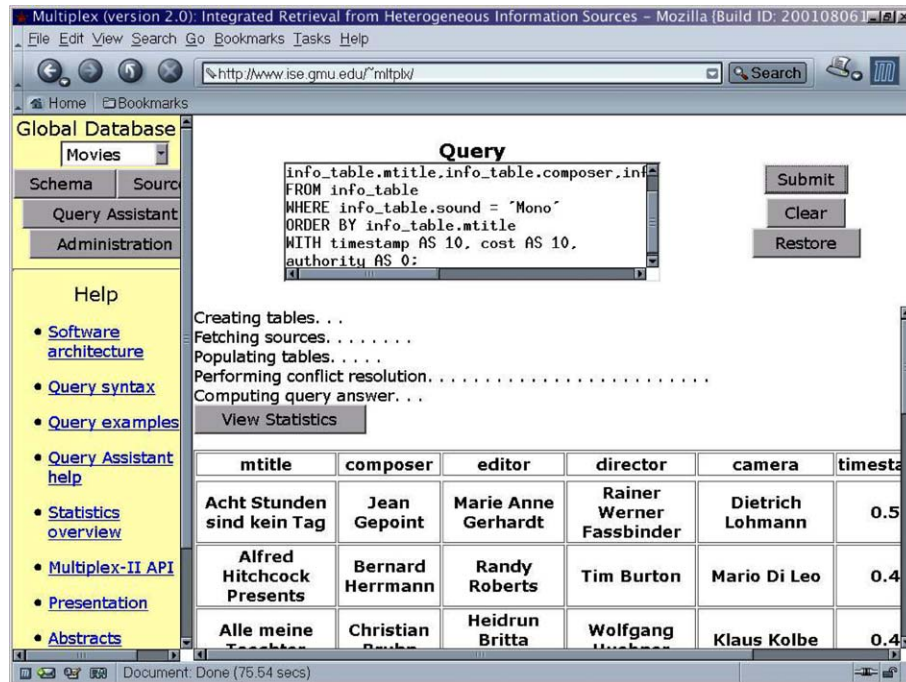


Fig. 5. Query result in Fusionplex.

of heterogeneity, these sources were stored in four different types of systems: a relational database, an object-oriented database, a plain file in a Unix system, and a Web-based resource available only through a form interface. The latter format is typical of how information is retrieved from Web information sources. These sources were fitted with simple “wrapping” software that implements a relational model “communication protocol” between the sources and the system. In one direction, each wrapper translates Fusionplex queries to the local query language; in the opposite direction, it assembles the source response in the tabular format that is understood by Fusionplex.

The overall architecture of Fusionplex and its tools provide for a flexible *integration service*. A remote user wishing to integrate several information sources (possibly, sources from this user’s own enterprise), logs into the server, provides it with the appropriate definitions, and can begin using its integration services right away. Future updates are fairly simple. For example, to integrate a new information source requires only a view definition (essentially, an SQL query) and a URL specification. Similarly, to expand the scope of the virtual database requires only a definition of a new virtual relation.

6.2. Experimentation

We describe in detail a large example that was tested in Fusionplex. The domain of the example is information on movies. Two Internet movie guides were chosen:

the Internet Movie Database (<http://us.imdb.com/>) and the All-Movie Guide (<http://allmovie.com/>). A virtual database called *Movies* was defined that integrates the information in these two sources with two relations (keys attributes are underlined):

Movie = (Title, Country, Year, Color, Production, Language, Runtime)

Credit = (Title, Sound, Genre, Director, Author, Composer, Camera, Editor)

Six contributions were extracted from the two sources and mapped into the global schema. All are join-selection-projection expressions over the two global relations. The global features are *timestamp*, *cost* and *authority*. The contributions are shown in Table 2.

A benchmark of 15 different queries was attempted against the virtual database. Three of these queries are shown in Table 3. The first query is a selection-projection of a single global relation; it involves four contributions. The second query demonstrates the use of aggregation on a global relation; it involves four contributions. The third query joins two global relations; it involves five contributions.

The number of tuples retrieved to process the queries in this benchmark ranged between 4065 and 7725, with an average of 6385 tuples per query. The number of tuples in the answers to these queries (the number of resolved tuples) ranged between 9 and 1956, with an average of 557 tuples per answer. The time required to process this benchmark was dominated completely by

Table 2
Contributions to the Virtual Database *Movies*

View Name: Ancient Movies
Features: (0.1,0.2,0.1)
Tuples: 90
select *Title, Country, Year, Color, Production, Language, Runtime, Sound, Genre, Director, Author, Composer, Camera*
from *Movie, Credit*
where *Movie.Title = Credit.Title and Year ≤ 1950*

View Name: Old Movies 1
Features: (0.3,0.1,0.9)
Tuples: 1,900
select *Title, Country, Year, Color, Production, Language, Runtime, Sound, Genre, Director, Author, Composer, Camera*
from *Movie, Credit*
where *Movie.Title = Credit.Title and Year > 1950 and Year ≤ 1980*

View Name: Old Movies 2
Features: (0.5,0.9,0.0)
Tuples: 2000
select *Title, Country, Year, Color, Production, Language, Runtime, Sound, Genre, Director, Author, Composer, Camera*
from *Movie, Credit*
where *Movie.Title = Credit.Title and Year > 1950 and Year ≤ 1980*

View Name: New Movies
Features: (0.4,0.4,0.4)
Tuples: 2260
select *Title, Country, Year, Color, Production, Language, Runtime, Sound, Genre, Director, Author, Composer, Camera*
from *Movie, Credit*
where *Movie.Title = Credit.Title and Year > 1980*

View Name: English Movies
Features: (0.5,1.0,1.0)
Tuples: 1400
select *Title, Country, Year, Color, Production, Language, Runtime, Sound, Genre, Director, Author, Composer, Camera*
from *Movie, Credit*
where *Movie.Title = Credit.Title and Language = "English"*

View Name: Short Movies
Features: (0.5,0.5,0.7)
Tuples: 75
select *Title, Country, Year, Color, Director*
from *Movie, Credit*
where *Movie.Title = Credit.Title and Runtime ≤ 60*

the time required to resolve conflicts. It ranged between 12s and 98s, with an average of 40s per query. Roughly speaking, the system resolved about 15 tuples per second. While this performance may not be sufficient for commercial applications, we believe that it affirms the feasibility of our fusion methods. It is expected that additional research on optimization and state-of-the-art equipment⁹ will improve performance considerably.

⁹ The equipment used in this experiment was a standard desktop computer with modest performance (e.g., processor speed: 1.2GHz, memory: 256MByte).

Table 3
Example Virtual Queries

Query 1:
select *Title, Year, Color, Runtime*
from *Movie*
where *Year ≥ 1970*
with *timestamp as 3, cost as 3, authority as 7*

Query 2:
select *Language, avg(Year), avg(Runtime)*
from *Movie*
where *Year ≥ 1960*
group by *Language*
order by *Language*
with *timestamp as 1, cost as 4, authority as 7*

Query 3:
select *Title, Country, Year, Genre, Director*
from *Movie, Credit*
where *Movie.Title = Credit.Title*
with *timestamp as 3, cost as 7, authority as 4*

7. Conclusion

7.1. Summary

In any realistic scenario of large-scale information integration, some of the information sources would be expected to “overlap” in their coverage. In such situations, it is practically unavoidable that they would occasionally provide inconsistent information. From users’ perspective, the desirable behavior of an integration system is to present them with answers that have been cleansed of all inconsistencies. This means that whenever multiple values contend for describing the same real world entity, they should be *fused* in a single value.

When humans are confronted with the need to choose a single value from a set of alternatives, they invariably consider the qualifications of the providers of these alternatives (and when they decide to take a simple average of the alternatives, or to choose one at random, it is usually because they conclude that the providers all have comparable qualifications).

The approach taken by the Fusionplex system formalizes these attitudes with the concept of source *features*, which quantify a variety of performance parameters of individual information sources.

Another observable behavior is that different individuals (or the same individuals in different situations) often apply different fusion policies. One individual might emphasize the importance of information recentness, whereas for another, cost might weight heaviest. As another example, in one situation, an individual would choose the average of the contending values, in another he would adopt the lowest, and in yet another he would pick the one that occurs most frequently.

Recognizing this, Fusionplex provides for powerful and flexible user control over the fusion process. In the two examples just mentioned, importance is

conveyed by means of feature weights, and the preferred resolution method is stated in a policy that combines quality thresholds, elimination guidelines and appropriate fusion functions.

The practicality of these principles has been demonstrated in a prototype implementation of Fusionplex. Besides its fusion strengths, this implementation provides a sound overall information integration environment, with support for information source heterogeneity, dynamic evolution of the information environment, quick ad-hoc integration, and intermittent source availability. These aspects are delivered in a client-server architecture that provides remote users with effective integration services, as well as convenient virtual database management tools.

7.2. Directions for further research

Research on Fusionplex is still continuing, with several issues currently under investigation, and other directions being considered. We discuss briefly nine such issues and directions.

Suppose fusions are implemented as linear combinations of the conflicting values (i.e., a generalization of *average*), and suppose users express their preferences in a *utility function* that is a linear combination of the features (as they do now). Not only can the conflicting values be ranked according to their utility to the user (this is done now in the procedure described in Section 5.4), it should also be possible to determine whether the utility of the fusion value is indeed higher than the utility of the existing values; i.e., if fusion is indeed profitable. Moreover, it should also be possible to generate automatically the *optimal* fusion policy: the fusion coefficients that optimize the utility function. Initial results on this issue are described in [20].

The resolution algorithm described in Section 5.4 assumes a *utility threshold* (denoted α), which must be provided. This threshold determines the tuples of sufficient quality that will participate in the inconsistency resolution process. It may be possible to choose this threshold at run-time, based on the data to be integrated, so that the overall utility of the result is *optimized*. Both of the last two directions remove the need for user input, by choosing parameters that optimize utility.

Inconsistency detection and fusion are performed at the attribute level. This implies that each of the values of a result tuple could come from a different version of the information. In some situations, this may be undesirable. Consider this example in which two different versions of postal address information exist, each with *City* and *Postal_code* attributes. Conceivably, the address generated by fusion could include a city from one version, and a postal code from another, resulting in incorrect information. In this example, *City* and *Postal_code* should constitute a single fusion unit. The meth-

odology should be extended to allow the specification of such units and to handle their fusion correctly.

During query translation, every contribution that is found to be relevant is materialized from its provider. Since providers are assumed to be able to deliver their contributions only as defined, inefficiencies may result. For example, a particular contribution may provide a large relation, of which the query might require only a few tuples. This is analogous to an Internet user downloading a large file, when in practice only a small portion of it is needed. The reason for this situation is that, in the interest of generality, we assumed that providers do not have capabilities for satisfying requests for *subsets* of their contributions. It would be useful to allow different classes of providers. Those that can only deliver their entire contributions, and those with capabilities of satisfying partial requests. Processing at the source would reduce considerably transmission and processing times for many queries.

As mentioned in Section 5.1, discrepancies might exist between the data “promised” in the view V of a contribution, and the data actually delivered when the URL is materialized. Fusionplex assumes that the information has been altered after the contribution has been plugged-in, and discards such contributions. It may be possible, however, to salvage some of the information by “repairing” contributions to correspond to their definitions.

Fusionplex extends standard SQL with new statements for expressing individual utility and resolution policies. It would be interesting to explore to what extent these tasks can be accomplished with standard SQL (e.g., **group by** and **having** constructs can be used to group, eliminate and fuse members of polytuples). The resultant simplicity of implementation might compensate in part for the loss of flexibility and control.

To participate in a Fusionplex virtual database, information providers must deliver their data in tabular format. XML [39] (the Extensible Markup Language) is quickly becoming a standard of data exchange. It would be beneficial to adopt XML as the communication protocol between Fusionplex and its information providers.

At times, users may benefit from having access to the entire set of alternative values after an inconsistency had been resolved. This would allow them to monitor the performance of their fusion (and possibly to redefine it). Presently, Fusionplex does not have this ability to “explain” its behavior.

Finally, a basic assumption in Fusionplex is that its contributing sources are associated with feature meta-data. Methodologies for collecting and maintaining these meta-data are outside the scope of this project, requiring separate research. But we mention here the challenge of feature maintenance. Feature values may need to be updated because of changes either in the actual content (e.g., *accuracy* or *timestamp*), or in external

parameters, either technological (e.g., *availability*) or commercial (e.g., *cost*). Arguably, updates of content-based features present the greater challenge. For these features, incremental methods should be devised that would *revise* them periodically or after updates, to avoid frequent recalculations.

References

- [1] S. Agarwal, A.M. Keller, G. Wiederhold, K. Saraswat, Flexible relation: an approach for integrating data from multiple, possibly inconsistent databases, in: Proceedings of ICDE-95, the Eleventh International Conference on Data Engineering, 1995, pp. 495–504.
- [2] P. Anokhin, Data inconsistency detection and resolution in the integration of heterogeneous information sources, Ph.D. thesis, School of Information Technology and Engineering, George Mason University, 2001.
- [3] P. Anokhin, A. Motro, Data integration: inconsistency detection and resolution based on source properties, in: Proceedings of FMII-01, International Workshop on Foundations of Models for Information Integration, September, 2001.
- [4] Y. Arens, C.A. Knoblock, W. Shen, Query reformulation for dynamic information integration, *Journal of Intelligent Information Systems* 6 (2/3) (1996) 99–130.
- [5] D. Barbara, H. Garcia-Molina, D. Porter, The management of probabilistic data, *IEEE Transactions on Knowledge and Data Engineering* 4 (5) (1992) 487–502.
- [6] E.F. Codd, Extending the database relational model to capture more meaning, *ACM Transactions of Database Systems* 4 (4) (1979) 397–434.
- [7] L.G. DeMichiel, Resolving database incompatibility: an approach to performing relational operations over mismatched domains, *IEEE Transactions on Knowledge and Data Engineering* 1 (4) (1989) 485–493.
- [8] O. Dunemann, I. Geist, R. Jesse, K.-U. Sattler, A. Stephanik, A database-supported workbench for information fusion: INFUSE, in: Proceedings of EDBT-02, 8th International Conference on Extending Database Technology, Lecture Notes in Computer Science, Springer-Verlag, 2002, pp. 756–758.
- [9] H. Galhardas, D. Florescu, E. Simon., C.-A. Saita, D. Shasha, Declarative data cleaning: language, model, and algorithms, in: Proceedings of VLDB-01, 27th International Conference on Very Large Data Bases, 2001, pp. 371–380.
- [10] H. Garcia-Molina, Y. Papakonstantinou, D. Quass, A. Rajaraman, Y. Sagiv, J.D. Ullman, V. Vassalos, J. Widom, The TSIMMIS approach to mediation: data models and languages, *Journal of Intelligent Information Systems* 8 (2) (1997) 117–132.
- [11] M.R. Genesereth, A.M. Keller, O. Duschka, Infomaster: an information integration system, in: Proceedings of ACM SIGMOD-97, International Conference on Management of Data, 1997, pp. 539–542.
- [12] M. Gertz, T.M. Özsu, G. Saake, K.-U. Sattler, Data Quality on the Web, Dagstuhl Seminar No. 03362, September 2003. <http://www.dagstuhl.de/03362/>.
- [13] Google. <http://www.google.com>.
- [14] A.Y. Halevy, Answering queries using views: a survey, *VLDB Journal* 10 (4) (2001) 270–294.
- [15] M.A. Hernandez, S.J. Stolfo, Real-world data is dirty: data cleansing and the merge/purge problem, *Data Mining and Knowledge Discovery* 2 (1) (1998) 9–37.
- [16] L. Kaufman, P.J. Rousseeuw, *Finding Groups in Data: an Introduction to Cluster Analysis*, John Wiley and Sons, 1990.
- [17] A.Y. Levy, A. Rajaraman, J.J. Ordille, Querying heterogeneous information sources using source descriptions, in: Proceedings of VLDB-96, 22nd International Conference on Very Large Data Bases, 1996, pp. 251–262.
- [18] E.-P. Lim, J. Srivastava, S. Shekhar, Resolving attribute incompatibility in database integration: an evidential reasoning approach, in: Proceedings of ICDE-94, 10th International Conference on Data Engineering, 1994, pp. 154–163.
- [19] A. Motro, Multiplex: a formal model for multidatabases and its implementation, in: Proceedings of NGITS-99, 4th International Workshop on Next Generation Information Technologies and Systems, Lecture Notes in Computer Science, vol. 1649, Springer-Verlag, 1999, pp. 138–158.
- [20] A. Motro, P. Anokhin, A.C. Acar, Utility-based resolution of data inconsistencies, in: Proceedings of IQIS-04, International Workshop on Information Quality in Information Systems, 2004, pp. 35–43.
- [21] A. Motro, I. Rakov, Estimating the quality of data in relational databases, in: Proceedings of the 1996 Conference on Information Quality, 1996, pp. 94–106.
- [22] A. Motro, I. Rakov, Not all answers are equally good: estimating the quality of database answers, in: T. Andreasen, H. Christiansen, H.L. Larsen (Eds.), *Flexible Query-Answering Systems*, Kluwer Academic Publishers, 1997, pp. 1–21.
- [23] F. Naumann, Quality-driven Query Answering for Integrated Information Systems, Lecture Notes in Computer Science, vol. 2261, Springer-Verlag, 2002.
- [24] F. Naumann, J. Freitag, M. Spilopoulou, Quality driven source selection using data envelope analysis, in: Proceedings of IQ-98, 3rd International Conference on Information Quality, 1998, pp. 137–152.
- [25] F. Naumann, M. Haeussler, Declarative data merging with conflict resolution, in: Proceedings of ICIQ-02, 7th International Conference on Information Quality, 2002, pp. 212–224.
- [26] F. Naumann, U. Leser, J.-C. Freytag, Quality-driven integration of heterogeneous information systems, in: Proceedings of the VLDB-99, 25th International Conference on Very Large Databases, 1999, pp. 447–458.
- [27] F. Naumann, C. Rolker, Assessment methods for information quality criteria, in: Proceedings of IQ-00, 5th International Conference on Information Quality, 2000, pp. 148–162.
- [28] F. Naumann, M. Scannapieco (Eds.), Proceedings of IQIS-04, International Workshop on Information Quality in Information Systems, 2004.
- [29] Y. Papakonstantinou, S. Abiteboul, H. Garcia-Molina, Object fusion in mediator systems, in: Proceedings of VLDB-96, 22nd International Conference on Very Large Data Bases, 1996, 413–424.
- [30] I. Rakov, Data quality and its use for resolving inconsistencies in multidatabase environments, Ph.D. thesis, School of Information Technology and Engineering, George Mason University, 1998.
- [31] R. Ramakrishnan, J. Gehrke, *Database Management Systems*, third Ed., McGraw-Hill, 2003.
- [32] E.M. Rasmussen, Clustering algorithms, in: W.B. Frakes, R.A. Baeza-Yates (Eds.), *Information Retrieval: Data Structures and Algorithms*, Prentice-Hall, 1992, pp. 419–442.
- [33] K.-U. Sattler, S. Conrad, G. Saake, Adding conflict resolution features to a query language for database federations, in: Proceedings of EFIS-00, the 3rd Workshop on Engineering Federated Information Systems, 2000, pp. 41–52.
- [34] V.S. Subrahmanian, S. Adali, A. Brink, R. Emery, J.J. Lu, A. Rajput, T.J. Rogers, R. Ross, C. Ward, HERMES: A heterogeneous reasoning and mediator system, 1994. <http://www.cs.umd.edu/projects/hermes/publications/abstracts/hermes.html>.
- [35] F.S.-C. Tseng, A.L.P. Chen, W.-P. Yang, A probabilistic approach to query processing in heterogeneous database systems, in: Proceedings of RIDE-TQP-92, Second International Workshop on Research Issues on Data Engineering: Transaction and Query Processing, 1992, pp. 176–183.

- [36] T.C. Redman, The impact of poor data quality in the typical enterprise, *Communications of the ACM* 41 (2) (1998) 78–82.
- [37] R. Wang et al. (Eds.), *Proceedings of the International Conference on Information Quality*, MIT Information Quality Program, 1996–2004.
- [38] R.Y. Wang, D. Strong, Beyond accuracy: what data quality means to data consumers, *Journal on Management of Information Systems* 12 (4) (1996) 5–34.
- [39] XML: Extensible Markup Language. <http://www.w3.org/XML>.