

Automatic Ontology Matching Using Application Semantics

Avigdor Gal, Giovanni Modica, Hasan Jamil, and Ami Eyal

■ We propose the use of application semantics to enhance the process of semantic reconciliation. Application semantics involves those elements of business reasoning that affect the way concepts are presented to users: their layout, and so on. In particular, we pursue in this article the notion of precedence, in which temporal constraints determine the order in which concepts are presented to the user. Existing matching algorithms use either syntactic means (such as term matching and domain matching) or model semantic means, the use of structural information that is provided by the specific data model to enhance the matching process. The novelty of our approach lies in proposing a class of matching techniques that takes advantage of ontological structures and application semantics. As an example, the use of precedence to reflect business rules has not been applied elsewhere, to the best of our knowledge. We have tested the process for a variety of web sites in domains such as car rentals and airline reservations, and we share our experiences with precedence and its limitations.

The ambiguous interpretation of concepts describing the meaning of data in data sources (for example, database schemata, extensible markup language [XML] document-type definitions [DTDs], Resource Description Framework [RDF] schemata, and hypertext markup language [HTML] form tags) is commonly known as *semantic heterogeneity*. Semantic heterogeneity, a well-known obstacle to data source integration, is resolved through a process of *semantic reconciliation*, which matches concepts from heterogeneous data sources. Traditionally, the complexity of semantic reconciliation required that it be performed by a human observer (a designer, a database admin-

istrator [DBA], or a user) (Hull 1997). However, manual reconciliation (with or without computer-aided tools) tends to be slow and inefficient in dynamic environments and, for obvious reasons, does not scale. Therefore, the introduction of the semantic web vision and the shift towards machine-understandable web resources has made clear the importance of automatic semantic reconciliation.

As an example, consider the web search, an information-seeking process conducted through an interactive interface. This interface may be as simple as a single input field (as in the case of a general-purpose search engine). Web interfaces may also be highly elaborate: consider a car rental or airline reservation interface containing multiple web pages, with numerous input fields, that are sometimes content dependent (for example, when a rented car is to be returned at the point of origin, no input field is required for the return location). A web search typically involves scanning and comparing web resources, either directly or by means of some information portal—a process hampered by their heterogeneity. Following the semantic web vision, semantic reconciliation should be inherent in the design of smart software agents for information seeking. Such agents can fill web forms and rewrite user queries by performing semantic reconciliation among different HTML forms.

To date, many algorithms have been proposed to support either semiautomatic or fully automatic matching of heterogeneous concepts in data sources. Existing matching algorithms make comparisons based on measures that are either syntactic in nature (such as term matching and domain matching) or based on model semantics. By model semantics, we

mean the use of structural information that is provided by the specific data model to enhance the matching process. For example, XML provides a hierarchical structure that can be exploited in identifying links among concepts and thus allow a smooth web search.

In this article, we propose the use of application semantics to enhance the process of semantic reconciliation. Application semantics involves those elements of business reasoning that affect the way concepts are presented to users, such as layout. In particular, we pursue in this article the notion of *precedence*, in which temporal constraints determine the order in which concepts are presented to the user.

All matching techniques aim at revealing latent semantics in data model descriptions and utilizing it to enhance semantic reconciliation. To illustrate the differences among syntactic measures and data model semantics on the one hand and application semantics on the other hand, consider a specific data model, XML, providing a domain description. Many matching techniques advocate the comparison of linguistic similarity, based on the assumption that within a single domain of discourse, terminology tends to be homogeneous. Linguistic similarity is based on terms that appear in the XML file. XML also has a hierarchical structure, allowing nesting of terms within other terms. This is a data model-specific feature (that does not exist in a relational model, for example), and may drive another approach towards matching. The underlying assumption here is that hierarchy is a feature designers of all applications can use to model the domain of discourse better and thus can be used to identify similarities.

We aim at moving beyond the data model, and to do so one has to analyze the domain of discourse (or several similar domains) to identify basic business rules and how they affect data modeling. As an example, say the XML file describes a car rental application. Analyzing this domain (and other similar domains, such as airline reservation systems) reveals temporal constraints that control the reservation process. For example, pickup location always precedes drop-off locations (both because renters typically drop off their rental at the same location and because the pickup location enforces constraints on the rest of the reservation, such as the availability of car types). Equipped with this observation, one can interpret the ordering within the XML file as a representative of such temporal constraints. To summarize, application semantics analysis starts at the application (and not at the data model as in the other approaches) and then is projected into the avail-

able data model to assist in the semantic reconciliation process.

The use of application semantics entails two immediate problems. First, it is likely that the data model does not support the application semantics features (or else they would have been used as data model semantics means). Therefore, there is the issue of formal representation of application semantics. Second, the lack of data model support means that algorithms that utilize application semantics are much harder to devise, having no underlying data model features upon which to be based.

To answer the first requirement of a rich data model for formal representation of application semantics, we choose to use ontologies. Ontologies are used as an interface conceptualization tool for representing model and application-level semantics to improve the quality of the matching process. Four ontological constructs are used in this work, namely terms, values, composition, and precedence. Terms, values, and composition are borrowed from Bunge (1977, 1979). Precedence, a unique feature of our model, represents the sequence in which terms are laid out within forms, imitating temporal constraints embedded in business rules.

In the general area of data integration, using a full-fledged ontology that is manually crafted to represent a domain of discourse with clear semantics and detached from a specific application is a rare privilege. More often than not, semantics is hidden in the application code, and only hints to it are divulged through interfaces and database schemata. Since our ontologies correspond directly to the semantics of the application, we propose (untraditionally) to abstract away ontologies from interfaces, thus exposing latent semantics. Therefore, composition can be extracted from the structure of a form, and precedence can be extracted from the ordering of elements in a form.

Given two ontologies (in the sense given above), algorithms to match terminologies in two web resources are needed. We propose syntactical comparison, based on terms and values, enhanced by basic information retrieval (IR) techniques for string matching. We also discuss what is needed to generate an algorithm that utilizes application semantics and discuss the difficulties in crafting such an algorithm, relating to the second problem presented above.

The novelty of our approach lies in the introduction of a sophisticated matching technique that takes advantage of ontological constructs and application semantics. In particular, the use of precedence to reflect business rules has not been applied elsewhere, to the best of our

knowledge. We have tested the process for a variety of web sites in domains such as car rentals and airline reservations and evaluated the performance of our algorithms. We highlight the benefits and limits of using the precedence construct as a guideline for future research into application semantics.

To support our research into application semantics, we developed *OntoBuilder*,¹ a tool that extracts ontologies from web applications and maps ontologies to answer user queries against data sources in the same domain. The input to the system is an HTML page representing the web site main page. Using *OntoBuilder*, HTML pages are parsed using a library for HTML/XML documents to identify form elements and their labels and to generate an ontology. Ontologies are then matched to produce a mapping using the algorithms presented in the “Ontology Matching” section. *OntoBuilder* supports an array of matching and filtering algorithms, and is extensible. It was developed using Java.

Research Background and Related Work

The study builds upon two existing bodies of research, namely heterogeneous databases and ontology design. Each is elaborated below.

Heterogeneous Databases

The evolution of organizational computing, from “islands of automation” into enterprise-level systems, has created the need to homogenize databases with heterogeneous schemata (referred to as *heterogeneous databases*). More than ever before, companies are seeking integrated data that go well beyond a single organizational unit. In addition, a high percentage of organizational data is supplied by external resources (for example, the web and extranets). Data integration is thus becoming increasingly important for decision support in enterprises. The growing importance of data integration also implies that databases with heterogeneous schemata face an ever-greater risk that their data integration process will not effectively manage semantic differences.

Current research into heterogeneous databases is largely geared toward manual (or semimanual semiautomatic at best) semantic resolution (such as Kahng and McLeod [1996] and Gal [1999]), which may not effectively scale in computational environments with dynamically changing schemata that require a rapid response. In addition, schema descriptions differ significantly among different domains. It is often said that the next great challenge in the se-

matic matching arena is the creation of a generalized set of automatic matching algorithms. Accordingly, the goal of this research is to propose the use of application semantics for automatic matching.

Over the past two decades, researchers in both academia and industry have advanced many ideas for reducing semantic mismatch problems, with the goal of lessening the need for manual intervention in the matching process. A useful classification of the various solutions proposed can be found in Rahm and Bernstein (2001). Of the categories presented there, we focus on those that deal with the algorithmic aspect of the problem.

The proposed solutions can be grouped into four main approaches. The first approach recommends adoption of information-retrieval techniques. Such techniques apply approximate, distance-based matching techniques, thus overcoming the inadequacy of exact, “keyword-based” matching. This approach is based on the presumption that attribute names can be mapped using similarity techniques. Attribute names are rarely, however, given in explicit forms that yield good matchings. Furthermore, they need to be complemented by either a lengthier textual description or an explicit thesaurus, which mandates greater human intervention in the process. *Protège* utilizes this method (among others) in the PROMPT algorithm, a semiautomatic matching algorithm that guides experts through ontology matching and alignment (Noy and Musen 2000).

A second approach involves the adoption of machine-learning algorithms that match attributes based on the similarity between their associated values. Most efforts in that direction (for example, Glue [Doan et al. 2002] and Autoplex [Berlin and Motro 2001]) adopt some form of a Bayesian classifier. In these cases, mappings are based on classifications with the greatest posterior probability, given data samples. Machine learning was recognized as playing an important role in reasoning about mappings in the work by Madhavan et al. (2002).

Third, several researchers have suggested the use of graph theory techniques to identify similarities among schemata, in which attributes are represented in the form of either a tree or a graph. To give but one example, the *TreeMatch* algorithm (Madhavan et al. 2002) utilizes XML DTD’s tree structure in evaluating the similarity of leaf nodes by estimating the similarity of their ancestors.

In a fourth approach, matching techniques from the first three groups are combined. Here, a weighted sum of the output of algorithms in these three categories is used to determine the

similarity of any two schema elements. Cupid (Madhavan, Bernstein, and Rahm 2001) and OntoBuilder are two models that support this hybrid approach. OntoBuilder, however, is the only framework, to the best of our knowledge, in which application semantics is used as a tool in matching heterogeneous schemata.

Ontology Design

The second body of literature we draw upon focuses on ontology design. An ontology is “a specification of a conceptualization” (Gruber 1993), in which conceptualization is an abstract view of the world represented as a set of objects. The term has been used in different research areas, including philosophy (where it was coined), artificial intelligence, information sciences, knowledge representation, object modeling, and most recently, e-commerce applications. For our purposes, an ontology can be described as a set of terms (vocabulary) associated with certain semantics and relationships. Typically, ontologies are represented using a description logic (Donini et al. 1996), in which subsumption typifies the semantic relationship between terms, or frame logic (Kifer, Lausen, and Wu 1995), in which a deductive inference system provides access to semistructured data.

The realm of information science has produced an extensive body of literature and practice in ontology construction (for example, Vickery [1966]). Other undertakings, such as the DOGMA project (Spyns, Meersman, and Jarrar 2002), provide an engineering approach to ontology management. Finally, researchers in the field of knowledge representation have studied ontology interoperability, resulting in systems such as Chimaera (McGuinness et al. 2000) and Protège (Noy and Musen 2000).

The body of research aimed at matching schemata by using ontologies has focused on interactive methods requiring human intervention, massive at times. In this work, we propose a fully automatic process that is a more scalable approach to semantic reconciliation. Our approach is based on analyzing model-dependent and application-level semantics to identify useful ontological constructs, followed by the design of algorithms to utilize these constructs in automatic schema matching. It is worth noting that automation carries with it a level of uncertainty as “the syntactic representation of schemas and data do not completely convey the semantics of different databases” (Miller, Haas, and Hernández 2000). In another paper (Gal et al. 2004), we have formally modeled the uncertainty inherent in an automatic semantic reconciliation and offered an evalua-

tion tool for the quality of algorithms that were designed for that purpose.

Ontological Constructs

The methodology for the process of schema matching is based on ontological analysis of application classes and the generation of appropriate ontological constructs that may assist in the matching process. We base the ontological analysis on the work of Bunge (1977, 1979). We adopt a conceptual modeling approach rather than a knowledge representation approach (in the AI sense). While the latter requires a complete reflection of the modeled reality for an unspecified intelligent task to be performed by a computerized system in the future (Borgida 1990), the former requires a minimal set of structures to perform a given task (a web search in this case). Therefore, we build ontologies from a given application (such as web forms) rather than with the assistance of a domain expert.

To exemplify the methodology, we focus on ontological constructs in the general task of the web search. We recognize the limited capabilities of HTML (and for that matter, also XML) in representing rich ontological constructs, and therefore we have eliminated many important constructs (for example, the class structure) simply because they cannot be realistically extracted from the content of web pages. Therefore, the ontological analysis of this class of applications yielded a subset of the ontological constructs provided by Bunge and added a new construct, which we term *precedence*, for posing temporal constraints.

Terms: We extract a set of terms² from a web page, each of which is associated with one or more form entries. Each form entry has a label that appears on the form interface and internal entry names that are not presented by the browser but are still available in HTML. The label provides the user with a description of the entry content. The latter is utilized for matching parameters in the data transfer process and therefore resembles the naming conventions for database schemata, including the use of abbreviations and acronyms. A term is a combination of both the label and the name. For example, *Airport Location Code* (`PICKUP_LOCATION_CODE`) is a term in the Avis reservation page, where *Airport Location Code* is the label and `PICKUP_LOCATION_CODE` is the entry name.

Values: Based on Bunge (1977), an attribute is a mapping of terms and value-sets into specific statements. Therefore, we can consider a combination of a term and its associated data entry (value) to be an attribute. In certain cases, the

value-set that is associated with a term is constrained using drop lists, check boxes, and radio buttons. For example, the entry labeled *Pick-Up Date* is associated with two value-sets: {Day, 1, 2, ..., 31} and {January, February, ..., December}. Clearly, the former is associated with the date of the month (and the value *Day* was added to ensure the user understands the meaning of this field) and the latter with the month (here, there is no need in adding a *Month* value, since the domain elements speak for themselves).

Composition: We differentiate atomic terms from composite terms. A composite term is composed of other terms (either atomic or composite). In the Avis reservation web page, all of the terms mentioned above are grouped under *Rental Pick-Up & Return Information*. It is worth noting that some of these terms are, in themselves, composite terms. For example, *Pick-Up Time* is a group of three entries, one for the hour, another for the minutes, and the third for either AM or PM.

Precedence: The last construct we consider is the precedence relationship among terms. In any interactive process, the order in which data are provided may be important. In particular, data given at an earlier stage may restrict the availability of options for a later entry. For example, the Avis web site determines which car groups are available for a given session using the information given regarding the pickup location and time. Therefore, once those entries are filled in, the information is sent back to the server and the next form is brought up. Such precedence relationships can usually be identified by the activation of a script, such as (but not limited to) the one associated with a SUBMIT button. It is worth noting that the precedence construct rarely appears as part of basic ontology constructs. This can be attributed to the view of ontologies as static entities whose existence is independent of temporal constraints. We anticipate that contemporary applications, such as the one presented in this article, will need to embed temporal reasoning in ontology construction.

The main difference between the first three constructs on the one hand, and the third construct on the other, is that the equivalence of the construct in the data model is given explicitly in the former but is only implicit in the latter. In our example, terms are explicitly available as labels and entry names, and values are explicitly available as value-sets. Composition is explicitly available in XML definitions through its hierarchical structure.³ The precedence construct, on the other hand, is only implicitly given, through the process of form submission.

It is worth noting that the recognition of useful ontological constructs is independent of the algorithms that are utilized to perform the reconciliation process. In the ensuing discussion, we shall demonstrate the usefulness of precedence in identifying correct mappings, yet discuss the difficulty of generating a good matching algorithm that avoids false positives and false negatives in the process.

Ontology Matching

In the matching process, a mapping is determined between two ontologies. To illustrate the complexity of the process, consider first the following example.

Example 1 (Ontology Matching)

Consider the Delta and American Airlines reservation systems (see figure 1). The left screen of figure 1 presents a form that contains two time fields, one for departure and the other for return. Due to bad design (or designer's error), the departure time entry is named *dept_time_1* while return time is named *dept_time_2*. Both terms carry an identical label, *Time*, since the context can be easily determined (by a human observer of course) from the positioning of the time entry with respect to the date entry. For the American Airlines reservation system (the right screen of figure 1), the two time fields of the latter were not labeled at all (counting on the proximity matching capabilities of an intelligent human being), and therefore were assigned, using composition by multiple term association, with the label *Departure Date* and *Return Date*. The fields were assigned the names *departureTime* and *returnTime*. Term matching would incur problems in differentiating the four terms (note that "dept" and "departure" do not match, either as words or as substrings).

We denote by web resource dictionary the set of terms extracted from a web resource (typically composed of several web pages within a single web site). Let $V = \{v_1, v_2, \dots, v_n\}$ and $U = \{u_1, u_2, \dots, u_m\}$ be two web resource dictionaries. The general matching process is conducted in two steps. First, pairwise matching yields a similarity measure for all pairs, and next a subset of the pairwise matching (dubbed a mapping) is selected as the "best" mapping between the two ontologies. Such a mapping may utilize some variation of a weighted bipartite graph matching (Galil 1986) if the required mapping is of a 1:1 nature. For a matching process that yields 1:*n* mappings, a simpler algorithm may be applied, in which a term in one dictionary is mapped into a term in another dictionary to which its similarity is maximized. Such an algorithm enables duplicate entries in one dictionary, yet does not allow the partition of a single value to several values.

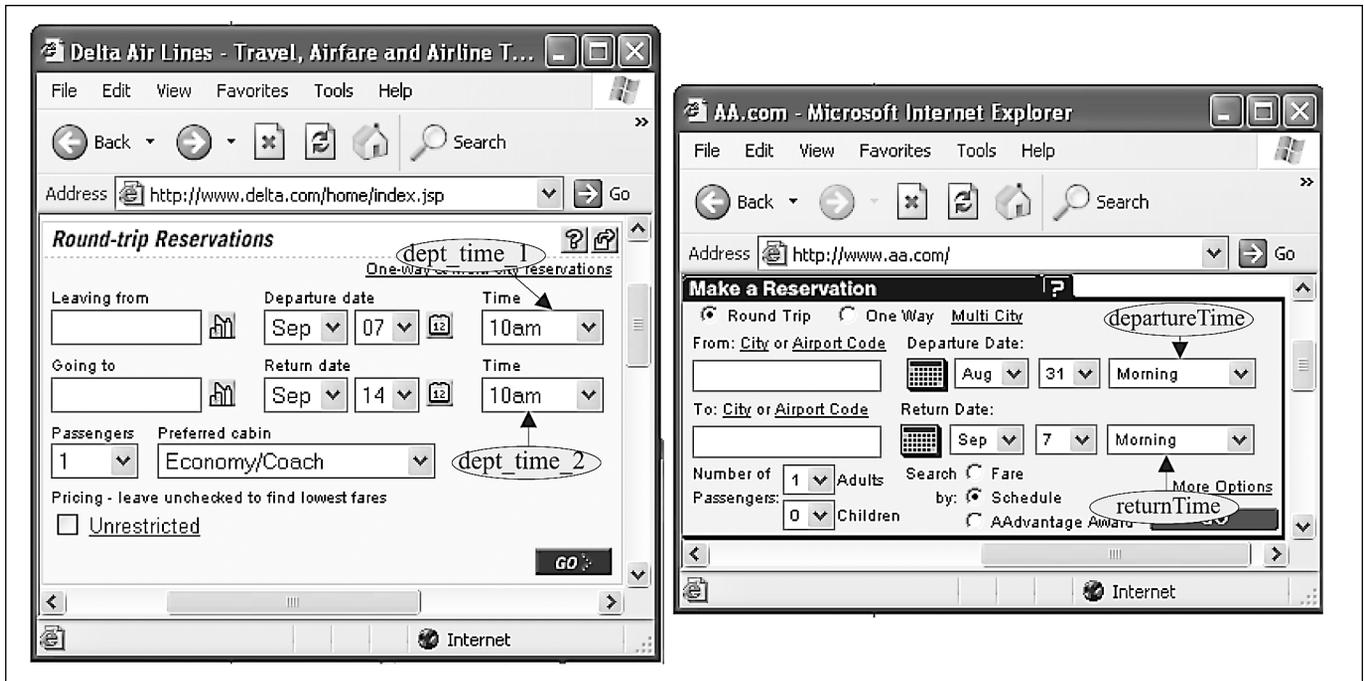


Figure 1. Delta Versus AA.

The process of ontology matching is formalized and discussed in depth in Gal et al. (2004). In particular, we have shown there that the specific methodology described herein is well suited to identifying the exact mapping (as perceived by a human observer) as the mapping with the highest sum (or average) of similarity measures of the selected term pairs.

The following sections focus on three methods for pairwise matching, namely term, value, and precedence matching. We omit the discussion of composition matching for the sake of brevity. A detailed algorithm is available in Modica (2002).

Syntactic Matching

In this section we present two syntactic methods for pairwise matching. We start the section with a discussion of term matching, then follow it with a discussion of value matching.

Term Matching

Term matching compares labels (verbal descriptions of a form entry) and names (the entry names as being sent to the server) to identify syntactically similar terms. To achieve better performance, terms are preprocessed using several techniques originating in IR research, including capitalization-based separation, ignorable character removal, dehyphenation, and stop-term removal.

We have applied two separate methods for term matching based on string comparison—word matching and string matching—as follows.

Word Matching. Two terms are matched and the number of common words is identified. The similarity of two terms t_1 and t_2 using word matching (dubbed $\mu(W, v, u)$) is defined as the ratio between the number of common words in t_1 and t_2 and the total number of unique words in terms t_1 and t_2 , providing a symmetric measure of the similarity of these two terms. The more common words the terms share, the more similar they are considered to be. For example, consider the terms $t_1 = \text{Pickup Location}$ and $t_2 = \text{Pick-up location code}$. The revised terms after preprocessing are $t_1 = \text{pickup location}$ and $t_2 = \text{pickup location code}$. The terms' similarity, using word matching, is computed as

$$\mu_{t_1, t_2}^w = \frac{2(\text{pickup, location})}{(\text{pickup, location, code})} = 66\%$$

Two words $w_1 \in t_1$ and $w_2 \in t_2$ are considered to be common if they are spelled the same, sound the same (soundex), or are considered synonyms, using a publicly available thesaurus such as WordNet.⁴ Mismatched terms can be presented to the user for manual matching. Every manual match identified by the user is accepted as a synonym and expands and enriches the thesaurus.

String Matching. We find the maximum common substring between two terms whose words have been concatenated by removing white

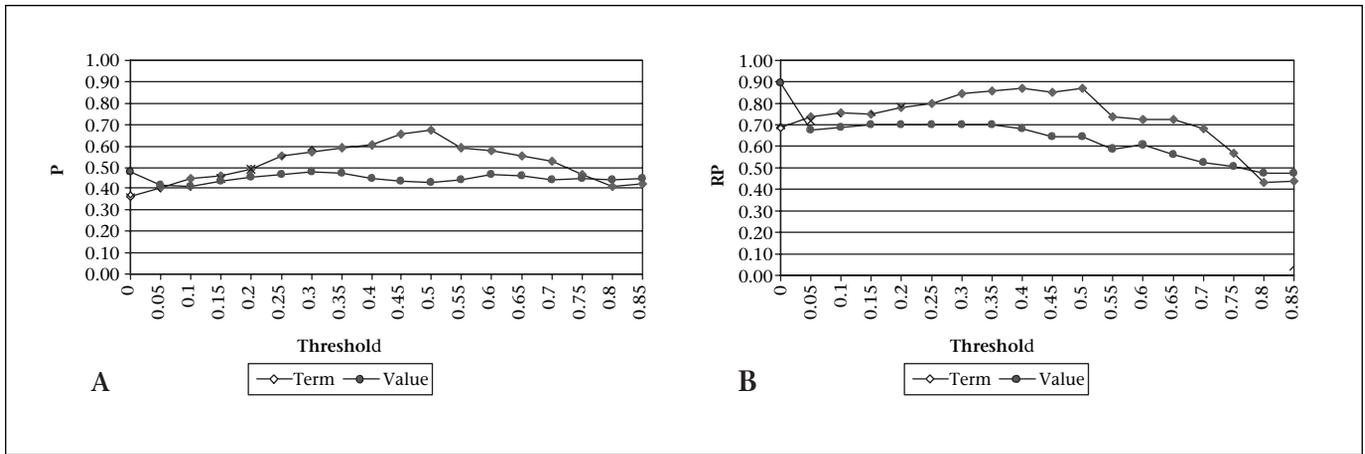


Figure 2. Precision and Relative Precision Versus Threshold.

spaces. The similarity of two terms using string matching (dubbed $\mu(S, v_i, u_j)$) is computed as the length of the maximum common substring as a percentage of the length of the longest of the two terms. As an example, consider the terms *airline information* and *flight airline info*, which after concatenating and removing white spaces become *airlineinformation* and *flightairlineinfo*, respectively. The maximum common substring is *airlineinfo*, and the effectiveness of the match is $\text{length}(\text{airlineinfo}) / \text{length}(\text{airlineinformation}) = 11/18 = 61$ percent.

We define a threshold (t^T) to identify a reasonable match. Any match with less than t^T is discarded. This threshold can be adjusted by the user.

For each pair, we compute four figures, two for labels, $\mu(W, L, v_i, u_j)$ and $\mu(S, L, v_i, u_j)$, and two for names, $\mu(W, N, v_i, u_j)$ and $\mu(S, N, v_i, u_j)$. We combine the figures into one figure, representing the strength of the match. Therefore, the similarity measure of a term v_i with a term u_j is computed as the weighted average

$$\begin{aligned} \mu_{v_i, u_j}^T &= \omega^{W,L} \mu_{v_i, u_j}^{W,L} + \omega^{S,L} \mu_{v_i, u_j}^{S,L} \\ &+ \omega^{W,N} \mu_{v_i, u_j}^{W,N} + \omega^{S,N} \mu_{v_i, u_j}^{S,N} \end{aligned} \quad (1)$$

where $\omega^{W,L}$, $\omega^{S,L}$, $\omega^{W,N}$, and $\omega^{S,N}$ are positive weights such that $\omega^{W,L} + \omega^{W,N} + \omega^{S,L} + \omega^{S,N} = 1$.

Experiments. We have conducted experiments to evaluate the performance of the term algorithm using two metrics, namely precision and relative precision. Let V and U be web resource dictionaries. U partitions V into two subsets V_1 and V_2 , such that V_1 is the set of all matchable terms and V_2 contains all those terms that cannot be matched with any term in U . Let M be a set of cardinality m , representing the set of all attributes in V that were matched by the algorithm. *Precision* (P) is the fraction of all found matches that are correct. It is computed as

$$P = \frac{|V_1 \cap M|}{m}$$

Relative precision is concerned with the ability of an algorithm to avoid false positives. Let t be a threshold. $V_1(t)$ is the set of all matchable terms among those terms for which the algorithm has given a similarity measure higher than t . *Relative precision* (RP) is computed to be

$$RP(t) = \frac{|V_1 \cap M(t)|}{|V_1(t)|}$$

The higher RP gets, the more efficient is the algorithm (at a given threshold) in avoiding false positives. It is worth noting that for $t = 0$, $V_1(t) = V_1$, and $RP(0)$ becomes recall, which is computed (using our terminology) as

$$\frac{|V_1 \cap M|}{|V_1|}$$

Figure 2 illustrates the performance of the term algorithm. Its performance varies from precision of 0.35 to 0.7. Its relative precision varies from 0.7 to 0.9. These results are good and can be attributed to the descriptive nature of labels in web forms. However, even at its peak, the term algorithm identifies 0.1 of the matches incorrectly. Example 1 has illustrated one such case that serves as a motivation to the presentation of the precedence construct. In Example 1, the term algorithm prefers matching both *Time(dept_time_1)* and *Time(dept_time_2)* of Delta with *Return Date(returnTime)* of American Airlines.

Value Matching

Value matching utilizes domain constraints to compute the similarity measure among terms. Whenever constrained value-sets are present, we can enhance our knowledge of the domain, since such constraints become valuable when comparing two terms that do not exactly

match through their labels. For example, the label corresponding to Avis's *Return Date* in Alamo's web site is *Dropoff Date*. The labels only partially match, and the words *Return* and *Dropoff* do not appear to be synonymic in general-purpose thesauri (*dropoff* is not even considered a word in English, according to the *Oxford English Dictionary*). Nevertheless, our matching algorithm matches these terms using their value-sets, since the term *Dropoff Date* has a value-set of $\{(Select), 1, 2, \dots, 31\}$ and the *Return Date* of Avis is associated with the value-set $\{Day, 1, 2, \dots, 31\}$.

It is our belief that designers would prefer constraining field domains as much as possible to minimize the effort of writing exception modules. Therefore, it is less likely (although known to happen occasionally) that a field with a drop-down list in one form will be designed as a text field in another form. In the case of a small-size domain, alternative designs may exist (for example, AM/PM may be represented as either a drop-down list or radio buttons). Since the extraction algorithm represents domains in a unified abstract manner, the end result is independent of the specific form of presentation.

Fields with select, radio, and check box options are processed using their value-sets. Therefore, different design methods act as no barrier in extracting the actual value sets. Value sets are preprocessed to result in generic domains. By recognizing separators in well-known data types, such as "/", "-", and "." in date structures, ":" in time structures, "(" in telephone numbers, "@" in e-mail addresses, and "http://" in URLs, domains can be partitioned into basic components, creating a compound term. The name of each new subterm is constructed as a concatenation of the existing name and the recognized domain type (for example, day). For example, the term *Pickup Date* (*pick_date*), which is recognized as a *date* field based on its domain entries, is further decomposed into three subterms: *Pickup Date* (*pick_date_day*), *Pickup Date* (*pick_date_month*), and *Pickup Date* (*pick_date_year*). It is worth noting that such preprocessing also affects term matching by generating additional terms and therefore is performed prior to term matching.

Similarity is calculated as the ratio between the number of common values in the two value sets and the total number of unique values in them. For example, suppose that $t_1 = \textit{Return time}$ and $t_2 = \textit{Dropoff time}$ with values $\{10:00AM, 10:30AM, 11:00AM\}$ and $\{10:00AM, 10:15AM, 10:30AM, 10:45AM, 11:00AM\}$, respectively. Preprocessing separates the domains into hour values ($\{10, 11\}$ versus $\{10, 11\}$), minutes values

($\{00, 30\}$ versus $\{00, 15, 30, 45\}$), and the value $\{AM\}$ (identical in both schemata). There is a perfect match in the hour domain, yet the minutes domains share two values (00 and 30) out of four (00, 15, 30, and 45). Thus, the similarity is calculated as $2/4 = 50$ percent. The power of value matching can be further highlighted using the case of *Dropoff Date* in Alamo and *Return Date* in Avis. These two terms have associated value sets $\{(Select), 1, 2, \dots, 31\}$ and $\{(Day), 1, 2, \dots, 31\}$ respectively, and thus their content-based similarity is $31/33 = 94$ percent, which improves significantly over their term similarity $(1(date) / 3(dropoff, date, return) = 33$ percent).

The domain recognition component can overcome differences of representation within the same domain. For example, we can apply transformations, such as converting a 24-hour representation into one of 12 hours. Thus, a domain $\{10:00, 11:00, 12:00, 13:00\}$ in a 24-hour representation can be transformed into three domains $\{1, 10, 11, 12\}$, $\{00\}$, and $\{AM, PM\}$ in a 12-hour representation.

Figure 2 illustrates the performance of the value algorithm, as a function of the threshold. The reasonable performance of the value algorithm is evident. What is not evident from this graph is that the value algorithm's performance varies much more than that of other algorithms. Clearly, for ontologies with many different data types, the value algorithm has good prediction capabilities (better than the term algorithm), while for ontologies in which many terms share the same domain, the value algorithm will find it much harder to predict correct mappings. The analysis of *relative precision* in figure 2b shows a repetition of the patterns in figure 2a. An interesting phenomenon is the ability of the value algorithm to outperform the term algorithm for a 0 threshold, with an average relative precision of 90 percent. This analysis can also serve in identifying optimal thresholds for various algorithms (in order to minimize false positives). Therefore, the value algorithm performs best at 0 threshold, while the term algorithm performs well in $[0.3, 0.5]$.

Returning to example 1, it is worth noting that value matching cannot differentiate the four possible combinations, since they share the same time domain. Therefore, other alternatives that better exploit the application semantics should be considered.

Precedence Matching

Let u_i and u_j be atomic terms in a web resource dictionary. If one of the following two conditions is satisfied, u_i precedes u_j : (1) u_i and u_j are associated with the same web page and u_i phys-

ically precedes u_j in the page; and (2) u_i and u_j are associated with two separate web pages, U_i and U_j , respectively, and U_i is presented to the user before U_j .

Evaluating the first condition is easily achieved when the page is extracted into a document object model (DOM) tree, a W3C standard that can be used in a fairly straightforward manner to identify form elements, labels, and input elements. The properties of the precedence relation are summarized in the following proposition.

Proposition 1

The precedence relation is irreflexive, antisymmetric, and transitive.

The precedence relationship, as presented in this article, serves as an estimation of the actual time constraints of a business process. For example, car rental companies would be likely to inquire about pickup information before return information. As yet another example, consider the advance search web pages of Lycos and Yahoo. The term algorithm has had difficulties in matching *member name (m_u)* with *yahoo i_d (login)*, giving it a score of 0.01. Instead, it preferred matching *member name (m_u)* with *list my new yahoo mail address free (mail directory)*, with a much higher score of 0.2. Precedence, on the other hand, indicates that login information precedes other terms in this category of web forms, putting it at the very beginning of the form.

Nevertheless, not all terms share precedence relationships. For example, there is no reason why either shipping address or invoice address should take precedence in a purchase order. To evaluate the difficulty of crafting a good matching algorithm, utilizing precedence, we have tested a simple algorithm using a technique we term *graph pivoting*. Given an atomic term v_i in a web resource dictionary V , we can compute the following two sets:

$$precede(v_i) = \{v_j \in V \mid v_j \text{ precedes } v_i\}$$

$$succeed(v_i) = \{v_j \in V \mid v_i \text{ precedes } v_j\}$$

It is worth noting that, following proposition 1, $precede(v_i) \cap succeed(v_i) = \emptyset$. Given two terms, v and u , from two web resource dictionaries V and U , respectively, we consider u and v to be pivots within their own ontologies. Therefore, we compute the similarity measure of matching $precede(v)$ with $precede(u)$, and $succeed(v)$ with $succeed(u)$. This computation is based on the syntactic similarity measures of the term and value algorithms. Presumably, terms will tend to match better if both those that precede them and those that succeed them do so. Our experiments show that the performance of this algorithm measures significantly lower in precision than the term algorithm (on-

ly 30–50 percent). The algorithm produces many false positive errors, suggesting that such an algorithm is put to better use in refuting possible matches than in supporting them.

Concluding Remarks

In this article, we have proposed the use of application semantics to enhance the process of ontology matching. Application semantics involves those elements of business reasoning that affect the way in which concepts are presented to users, for example through their layout. In particular, we have introduced the precedence ontological construct, in which temporal constraints determine the sequence of concepts presented to the user. While the article has suggested the extraction of ontologies from HTML forms, we consider the use of ontologies to be essential for the broad area of web search. Current search engines (in particular Google) have applied IR techniques in matching documents with user queries. We believe that the addition of structures such as precedence to search engines, whenever suitable, would enhance the precision of the search process. We leave this as an open research question. In particular, we will explore the use of additional ontology structures to improve the effectiveness of the matching process.

It is our conjecture that using application semantics as a means for semantic reconciliation can be generalized beyond its application to HTML web forms. For example, the relational model has little ability to represent application semantic means such as precedence. However, many relational databases are interfaced nowadays through the use of HTML forms, for which precedence (and other application semantics) can increase the success of semantic reconciliation. Also, analysis of typical queries for a given application reveals information regarding the typical use of concepts, which can be further utilized in the semantic reconciliation process. We plan on investigating the methods illustrated above in future research.

While precedence has proven itself in certain instances, a good algorithm is still needed to extract this knowledge and put it to use, as our experiments show. The conceptual framework we provide, however, opens the door to more application-semantic concepts to be introduced and used in the ontology matching process.

We aim at continually improving the proposed algorithms. For example, the use of a linear algorithm for finding the maximal substrings and superstrings of two given strings was suggested in the context of bioinformatics

(Ruzzo and Tompa 1999). Embedding a variation of this algorithm in our system may reduce the complexity of string matching. Finally, we intend to research in depth the problem of complex query rewriting in a heterogeneous schemata setting, using data-type identification and domain normalization. The method proposed in this work serves as a promising starting point, yet a more thorough methodology is yet to be developed.

Research complementing the present article provides sufficient conditions for matching algorithms to identify exact mappings, as conceived by an expert. This work is reported in Anaby-Tavor, Gal, and Trombetta (2003) and Gal et al. (2004).

Acknowledgements

Our thanks is given to Louiqa Raschid and An-Hai Doan for useful discussions. This research was partially supported by the Fund for the Promotion of Research at the Technion (191-496) and by the Fund of the Vice President for Research at the Technion (191-507). Giovanni Modica and Hasan Jamil's research was partially supported by National Science Foundation EPSCoR Grants (EPS0082979 and EPS-0132618), a USDA-ARS Cooperative Agreement grant (CRIS-6406-21220-005-15S), and a Southwest Mississippi Resource Conservation and Development Grant (01050412). Avigdor Gal's research was partially supported by the IBM Faculty Award (2002). We thank Ido Peled, Haggai Roitman, and the class of "Information Systems and Knowledge Engineering Seminar," fall semester, 2002, for their assistance in collecting and analyzing the data.

Notes

1. Available at <http://ie.technion.ac.il/OntoBuilder>.
2. The choice of words to describe ontological constructs in Bunge's work had to be general enough to cover any application. We feel that the use of *thing*, which may be reasonable in a general framework, can be misleading in this context. Therefore, we have decided to replace it with the more concrete description of *term*.
3. Forms are given in HTML, which does not have a composition construct per se. Yet, our methodology transforms the HTML code into an XML definition, to be utilized in the reconciliation process.
4. <http://www.cogsci.princeton.edu/~wn/>.

References

Anaby-Tavor, A.; Gal, A.; and Trombetta, A. 2003. Evaluating Matching Algorithms: The Monotonicity Principle. Paper presented at the IJCAI-03 Workshop on Information Integration on the Web, Acapulco, Mexico, August 9–10 (www.isi.edu/info-agents/workshops/ijcai03/proceedings.htm).

- Berlin, J.; and Motro, A. 2001. Autoplex: Automated Discovery of Content for Virtual Databases. In *Cooperative Information Systems: Ninth International Conference, CoopIS 2001*, volume 2172 of Lecture Notes in Computer Science, 108–122. Berlin: Springer-Verlag.
- Borgida, A. 1990. Knowledge Representation, Semantic Data Modeling: What's the Difference? In Proceedings of the Ninth International Conference on Entity-Relationship Approach (ER'90), 1–2. Lausanne, Switzerland: ER Institute.
- Bunge, M. 1977. *Treatise on Basic Philosophy: Ontology I: The Furniture of the World*. Volume 3. New York: D. Reidel Publishing.
- Bunge, M. 1979. *Treatise on Basic Philosophy: Ontology II: A World of Systems*. Volume 4. New York: D. Reidel Publishing.
- Doan, A.; Madhavan, J.; Domingos, P.; and Halevy, A. 2002. Learning to Map between Ontologies on the Semantic Web. In Proceedings of the Eleventh International Conference on the World Wide Web, 662–673. New York: Association for Computing Machinery.
- Donini, F. M.; Lenzerini, M.; Nardi, D.; and Schaerf, A. 1996. Reasoning in Description Logic. In Principles on Knowledge Representation, Studies in Logic, Languages and Information, ed. G. Brewka, 193–238. Stanford, CA: CSLI Publications.
- Gal, A. 1999. Semantic Interoperability in Information Services: Experiencing with CoopWARE. *SIGMOD Record* 28(1): 68–75.
- Gal, A.; Anaby-Tavor, A.; Trombetta, A.; and Montesi, D. 2004. A Framework for Modeling and Evaluating Automatic Semantic Reconciliation. *VLDB Journal* 13(4).
- Galil, Z. 1986. Efficient Algorithms for Finding Maximum Matching in Graphs. *ACM Computing Surveys* 18(1)(March): 23–38.
- Gruber, T. R. 1993. A Translation Approach to Portable Ontology Specifications. *Knowledge Acquisition* 5(2): 199–220.
- Hull, R. 1997. Managing Semantic Heterogeneity in Databases: A Theoretical Perspective. In Proceedings of the ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS), 51–61. New York: Association for Computing Machinery.
- Kahng, J.; and McLeod, D. 1996. Dynamic Classification Ontologies for Discovery in Cooperative Federated Databases. In Proceedings of the First International Foundation on Cooperative Information Systems (IFCIS) International Conference on Cooperative Information Systems (CoopIS'96), 26–35. Brussels, Belgium: IFCIS.
- Kifer, M.; Lausen, G.; and Wu, J. 1995. Logical Foundation of Object-Oriented and Frame-Based Languages. *Journal of the ACM* 42(4): 741–843.
- Madhavan, J.; Bernstein, P. A.; Domingos, P.; and Halevy, A. Y. 2002. Representing and Reasoning about Mappings between Domain Models. In Proceedings of the Eighteenth National Conference on Artificial Intelligence and Fourteenth Conference on Innovative Applications of Artificial Intelligence (AAAI/IAAI), 80–86. Menlo Park, CA: AAAI Press.
- Madhavan, J.; Bernstein, P.A.; and Rahm, E. 2001.

Generic Schema Matching with Cupid. In *Proceedings of the International Conference on Very Large Data Bases (VLDB)*, 49–58. San Francisco: Morgan Kaufmann Publishers.

McGuinness, D. L.; Fikes, R.; Rice, J.; and S. Wilder, S. 2000. An Environment for Merging and Testing Large Ontologies. In *Proceedings of the Seventh International Conference on Principles of Knowledge Representation and Reasoning (KR2000)*. San Francisco: Morgan Kaufmann Publishers.

Miller, R. J.; Haas, L. M.; and Hernández, M. A. 2000. Schema Mapping as Query Discovery. In *Proceedings of the International Conference on Very Large Data Bases (VLDB)*, ed. A. El Abbadi, M. L. Brodie, S. Chakravarthy, U. Dayal, N. Kamel, G. Schlageter, and K.-Y. Whang, 77–88. San Francisco: Morgan Kaufmann Publishers.

Modica, G. 2002. A Framework for Automatic Ontology Generation from Autonomous Web Applications. Master's thesis, Mississippi State University, Mississippi State, MS.

Noy, N. F.; and Musen, M. A.. 2000. PROMPT: Algorithm and Tool for Automated Ontology Merging and Alignment. In *Proceedings of the Seventeenth National Conference on Artificial Intelligence (AAAI-2000)*, 450–455. Menlo Park, CA: AAAI Press.

Rahm, E. and Bernstein, P. A. 2001. A Survey of Approaches to Automatic Schema Matching. *VLDB Journal* 10(4): 334–350, 2001.

Ruzzo, W. L.; and Tompa, M. 1999. A Linear Time Algorithm for Finding All Maximal Scoring Subsequences. In *Proceedings of the Seventh International Conference on Intelligent Systems for Molecular Biology*, ed. T. Lengauer, R. Schneider, P. Bork, D. L. Brutlag, J. I. Glasgow, H.-W. Mewes, and R. Zimmer, 234–241. Menlo Park, CA: AAAI Press.

Spyns, P.; Meersman, R.; and Jarrar, M. 2002. Data Modeling Versus Ontology Engineering. *SIGMOD Record*, 31(4) 12–17.

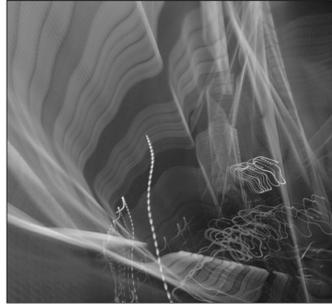
Vickery, B. C. 1966. Faceted Classification Schemes. Technical Report, Graduate School of Library Service, Rutgers, The State University, New Brunswick, NJ.



Avigdor Gal is a senior lecturer at the Technion—Israel Institute of Technology. He obtained his D.Sc. at the Technion. The focus of his work is on data integration and schema matching in databases and web environment. His e-mail address is avigal@ie.technion.ac.il.



Ami Eyal is a master's degree candidate at the Faculty of Industrial Engineering and Management, Technion—Israel Institute of Technology. He earned his B.Sc. in 1999 in industrial engineering and management from the Technion. Eyal's research is in the area of knowledge and information systems en-



**Journal of Artificial
Intelligence Research**

Volume 22, 2004
Martha E. Pollack
Editor-in-Chief
Moshe Tennenholz
Associate Editor-in-Chief

Now available from AAAI Press!
www.aaai.org/Press/Journals/JAIR

gineering. Research interests include effective methods for data integration and schema matching. His e-mail address is eyalami@tx.technion.ac.il.



Giovanni Modica is a computer scientist with experience both in academics and industry. His main areas of interest are databases, data integration, and business intelligence. For the last couple of years he has specialized in CRM systems, working as a team leader and developer for projects in different industries.

Modica has an M.Sc. degree from the Computer Science Department at Mississippi State University. His e-mail address is modicag@hotmail.com.



Hasan Jamil is a member of the faculty in the Department of Computer Science, Wayne State University. He earned his Ph.D. degree in computer science from Concordia University, Canada, and his M.S. and B.S. degrees in applied physics and electronics from the University of Dhaka, Bangladesh. He was also a member of the computer science faculty at Concordia University, Macquarie University, and Mississippi State University before joining Wayne State University in 2003. His current research interests are in the areas of databases, bioinformatics, and knowledge representation. He is currently a member of the editorial board of the *ACM Applied Computing Review*, and the chair of the IFIP TC 5 Bioinformatics Special Interest Group. He can be reached at jamil@acm.org.



5th International Symposium on Smart Graphics

August 22-24, 2005, near Munich, Germany

The fifth International Symposium on Smart Graphics will bring together researchers from computer graphics, visualization, art & graphics design, cognitive psychology and artificial intelligence, all working on different aspects of computer-generated graphics. This year's meeting will be held in the beautifully calm and serene atmosphere of the Frauenwoerth cloister near Munich, Germany.

Advances and breakthroughs in computer graphics have made visual media the basis of the modern user interface, and it is clear that graphics will play a dominant role in the way people communicate and interact with computers in the future. Indeed, as computers become more and more pervasive, and display sizes both increase and decrease, new and challenging problems arise for the effective use and generation of computer graphics.

Recent advances in this field have allowed AI researchers to integrate graphics in their systems, and on the other hand, many AI techniques have matured to the point of being easily used by non specialists. These very techniques are likely to be the vehicle by which both principles from graphics design and the results of research into cognitive aspects of visual representations will be integrated in next generation graphical interfaces.

Important Dates:

April 24 Submission deadline
May 23 Notification of review results
May 31 Camera ready copy due
Aug 22-24 Smart Graphics Symposium

Organizing Committee:

Andreas Butz (University of Munich)
Brian Fisher (Univ. of British Columbia)
Antonio Krueger (University of Muenster)
Patrick Olivier (University of Newcastle)

Proceedings:

Published as Springer LNCS

Hosted by:

University of Munich

Symposium Venue:

Frauenwoerth Cloister on Frauenchiemsee
island near Munich, Germany

In cooperation with:

ACM (pending), AAAI, Eurographics
Association (EG)

Symposium website:

www.smartgraphics.org

Help us celebrate AAAI's Twenty-Fifth Anniversary!

*Join us at AAAI-05
in Pittsburgh
July 9-13, 2005*

www.aaai.org/Conferences/National/2005/