

Constructing Data Graphs for Keyword Search* †

Konstantin Golenberg¹ and Yehoshua Sagiv²

¹ The Hebrew University, Jerusalem 91094, Israel,
konstg01@cs.huji.ac.il

² The Hebrew University, Jerusalem 91094, Israel,
sagiv@cs.huji.ac.il

Abstract. A data graph is a convenient paradigm for supporting keyword search that takes into account available semantic structure and not just textual relevance. However, the problem of constructing data graphs that facilitate both efficiency and effectiveness of the underlying system has hardly been addressed. A conceptual model for this task is proposed. Principles for constructing good data graphs are explained. Transformations for generating data graphs from RDB and XML are developed. The results obtained from these transformations are analyzed. It is shown that XML is a better starting point for getting a good data graph.

Keywords: Data-graph construction, keyword search, RDB, XML

1 Introduction

Considerable research has been done on effective algorithms for keyword search over data graphs (e.g., [3, 4, 7, 10–12, 14, 17]). Usually, a data graph is obtained from RDB, XML or RDF by a rather simplistic transformation. In the case of RDB [3, 6, 12], tuples are nodes and foreign keys are edges. When the source is XML [11, 13], elements are nodes, and the edges reflect the document hierarchy and IDREF(S) attributes.

In many cases, the source data suffers from certain anomalies and some papers (e.g., [13, 15]) take necessary steps to fix those problems. For example, when citations are represented by XML elements, they should be converted to IDREF(S) attributes. As another example, instead of repeating the details of an author in each one of her papers, there should be a single element representing all the information about that author and all of her papers should reference that element. These are examples of necessary transformations on the source data. If they are not done, existing algorithms for keyword search over data graphs will not be able to generate meaningful answers.

Once a source data is ameliorated, it should be transformed into a graph. The literature hardly discusses how it should be done. In [3, 14], the source is an RDB and the naive approach mentioned earlier is used (i.e., tuples are nodes and foreign keys are edges). In [11, 20], the source data is XML and the simplistic transformation described at the beginning of this section is applied. In [2, 5, 12, 16, 18], they do not mention any

* This work was supported by the Israel Science Foundation (Grant No. 1632/12).

† This paper is the full version of [9]. The final publication is available at Springer via http://dx.doi.org/10.1007/978-3-319-44406-2_33.

details about the construction of data graphs. The lack of a thoughtful discussion in any of those papers is rather surprising, because the actual details of constructing a data graph have a profound effect on both the efficiency and the quality of keyword search, regardless of the specific algorithms and techniques that are used for generating answers and ranking them.

Construction of effective data graphs is not a simple task, since the following considerations should be taken into account. For efficiency, a data graph should be as small as possible. It does not matter much if nodes have large textual contents, but the number of nodes and edges is an important factor. However, lumping together various entities into a single node is not a good strategy for increasing efficiency, because answers to queries would lose their coherence.

The structure of a data graph should reflect succinctly the semantics of the data, or else answers (which are subtrees) would tend to be large, implying that finding them would take longer and grasping their meaning quickly would not be easy.

An effective engine for keyword search over data graphs must also use information-retrieval techniques. Those tend to perform better on large chunks of text, which is another reason against nodes with little content.

In this paper, we address the problem of how to construct data graphs in light of the above considerations. In Section 4, we develop transformations for constructing data graphs from RDB and XML. In Section 5, we show that the format of the source data (i.e., RDB or XML) has a significant impact on the quality of the generated data graph. Moreover, XML is a better starting point than RDB. This is somewhat surprising given the extensive research that was done on designing relational database schemes.

As a conceptual guideline for constructing a good data graph, we use the OCP model [1], which was developed for supporting a graphical display of answers so that their meaning is easily understood. In Section 3, we explain why the OCP model is also useful as a general-purpose basis for constructing data graphs in a way that takes into account all the issues mentioned earlier.

In summary, our contributions are as follows. First, we enunciate the principles that should guide the construction of data graphs. Second, we develop transformations for doing so when the source data is RDB or XML. These transformations are more elaborate than the simplistic approach that is usually applied. Third, we show how the format of the source data impacts the quality of the generated graphs. Moreover, we explain why XML is a better starting point than RDB.

Our contributions are valid independently of a wide range of issues that are not addressed in this paper, such as the algorithm for generating answers and the method for ranking them. We only assume that an answer is a non-redundant subtree that includes all the keywords of the query. However, our results still hold even if answers are subgraphs, as sometimes done.

A presentation that gives motivation for the work of this paper is given in [9].

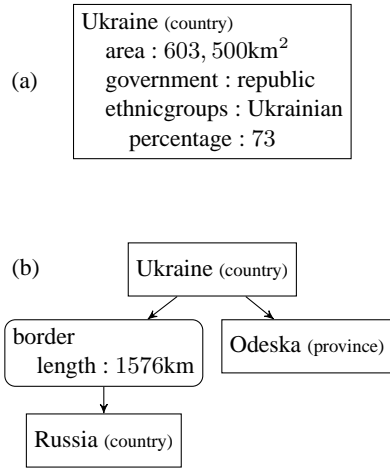


Fig. 1. An object and a tiny snippet of a data graph (not all properties are shown)

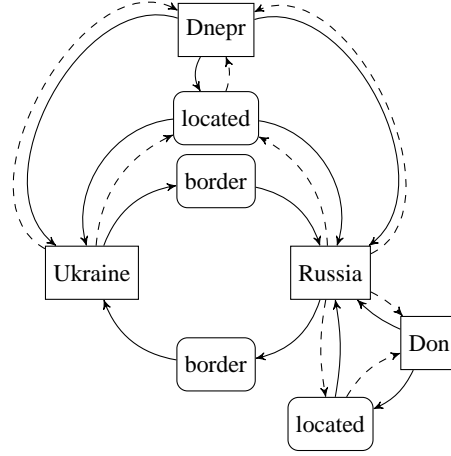


Fig. 2. A tiny portion of Mondial

2 Preliminaries

2.1 The OCP Model

The *object-connector-property* (OCP) model for data graphs was developed in [1] to facilitate an effective GUI for presenting subtrees. (As explained in the next section, those subtrees are answers to keyword search over data graphs.) In the OCP model, objects are entities and connectors are relationships. We distinguish between two kinds of connectors: *explicit* and *implicit*. Objects and explicit connectors can have any number of properties. Two special properties are *type* and *name*.

Parts (a) and (b) of Figure 1 show an object and a snippet of a data graph, respectively. An object is depicted as a rectangle with straight corners. The top line of the rectangle shows the name and type of the object. The former appears first (e.g., *Ukraine*) and the latter is inside parentheses (e.g., *country*). The other properties appear as pairs consisting of the property's name and value, as shown in Figure 1(a). Observe that properties can be nested; for example, the property *percentage* is nested inside *ethnicgroup*. Nesting is indicated in the figure by indentation.

An implicit connector is shown as a directed edge between two objects. Its meaning should be clear from the context. In Figure 1(b), the implicit connector from *Ukraine* to *Odeska* means that the latter is a province in the former.

An explicit connector is depicted as a rectangle with rounded corners. It has at most one incoming edge from an object and any positive number of outgoing edges to some objects. An explicit connector has a type, but no name, and may also possess other properties. Figure 1(b) shows an explicit connector of type *border* from *Ukraine* to *Russia* that has the property *length* whose value is 1576km.

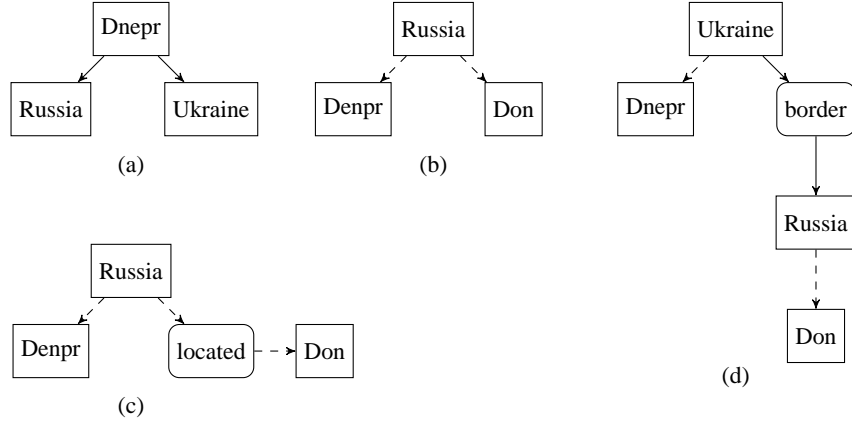


Fig. 3. Answers to queries

2.2 Answers to Keyword Search

We consider keyword search over a directed data graph G . (A data graph must be directed, because relationships among entities are not always symmetric.) A *directed subtree* t of G has a unique node r , called the *root*, such that there is exactly one directed path from r to each node of t .

A query Q over a data graph G is a set of keywords, namely, $Q = \{k_1, \dots, k_n\}$. An *answer* to Q is a directed subtree t of G that contains all the keywords of Q and is nonredundant, in the sense that no proper subtree of t also contains all of them.

For example, consider Figure 2, which shows a snippet of the data graph created from the XML version of the Mondial dataset,³ according to the transformation of Section 4.2. To save space, only the name (but not the type) of each object is shown. The dashed edges should be ignored for the moment. The subtree in Figure 3(a) is an answer to the query $\{\text{Dnepr}, \text{Russia}, \text{Ukraine}\}$. There are additional answers to this query, but all of them have more than three nodes and at least one explicit connector.

For the query $\{\text{Dnepr}, \text{Don}\}$, there is no answer (with only solid edges) saying that Dnepr and Don are rivers in Russia, although the data graph stores this fact. The reason is that the connectors (in the data graph of Figure 2) have a symmetric semantics, but the solid edges representing them are in only one direction. The only exception is the connector *border*, which is already built into the graph in both directions (between Russia and Ukraine). In order not to miss answers, we add *opposite edges* when symmetric connectors do not already exist in both directions. Those are shown as dashed arrows. Now, there are quite a few answers to the query $\{\text{Dnepr}, \text{Don}\}$ and Figure 3(b)–(d) shows three of them. The first two of those say that Dnepr and Don are rivers in Russia. These two answers have the same meaning, because the relationship between a river and a country is represented twice: by an implicit connector and by the explicit connector *located*. The answer in Figure 3(d) has a different meaning,

³ <http://www.dbis.informatik.uni-goettingen.de/Mondial/>

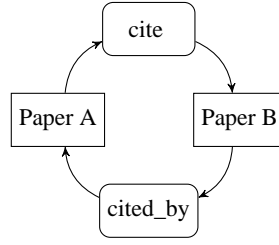


Fig. 4. An asymmetric connector and its inverse



Fig. 5. A single connector node for border

namely, Dnepr and Don are rivers in Ukraine and Russia, respectively, and there is a border between these two countries.

To generate relevant answers early on, weights are assigned to the nodes and edges of a data graph. Existing algorithms (e.g., [3, 7, 8, 10] enumerate answers in an order that is likely to be correlated with the desired one. Developing an effective weighting scheme is highly important, but beyond the scope of this paper.

2.3 Why Data Graphs are Directed

Data graphs must be directed because some relationships are asymmetric. For example, to represent citations among papers, we need two different types of connectors, as shown in Figure 4. In contrast, one connector type is sufficient for representing borders.

When a relationship is symmetric, it is redundant to use a different connector node for each direction, which is the case with `border` in Figure 2. It is better to represent a border between two countries as in Figure 5.

Over the data graph of Figure 5, the following are exactly two answers to the query `{Russia, Ukraine}`.

- `Ukraine → border → Russia`
- `Russia --> border --> Ukraine`

As directed subtrees, these answers are distinct. However, they carry the same information. Hence, we eliminate duplicates (similarly to [3]) by treating an answer as a set of undirected edges. That is, two answers are the same if they have the same set of undirected edges. Equality of undirected edges is determined as follows. Each node has a unique id (which is internal to the system). Thus, two edges are identical if they are the same unordered pair of id's.

Over the data graph of Figure 2, however, there are two distinct `border` nodes between Ukraine and Russia. Hence, the following two answers

- `Ukraine → border → Russia`
- `Russia → border → Ukraine`

are distinct even when viewed as undirected subtrees. To eliminate duplicates also in this case, we need to consider two connector nodes as equal if they have the same type, rather than the same id.

Even when a connector type is asymmetric, it is redundant to present both directions. For example, given the data graph of Figure 4, the following two subtrees carry the same information, in spite of having nodes of different types.

- Paper A \rightarrow cite \rightarrow Paper B
- Paper B \rightarrow cited_by \rightarrow Paper A

To eliminate one of these two as a duplicate, we need to treat two connector nodes as equal if one has the inverse type of the other.

3 Advantages of the OCP Model

In this section, we discuss some of the advantages of the OCP model. In a naive approach of building a data graph, there is only one type of nodes (i.e., no distinction between objects and connectors). Moreover, sometimes there is even a separate node for each property. This approach suffers from three drawbacks. First, from the implementation's point of view, this is inefficient in both time and space. That is, even if there is not much data, the number of nodes and edges is likely to be large. As a result, searching a data graph for answers would take longer (than the alternative described later in this section). In addition, if all the processing is done in main memory, the size of the data graph is more likely to become a limiting factor.

The second drawback of the naive approach is from the user's point of view. A meaningful answer is likely to have quite a few nodes; hence, displaying it graphically in an easily understood manner is rather hard. Another problem is the following. The definition of an answer is intended to avoid redundant parts in order to cut down the search space. However, sometimes an answer must be augmented to make it clear to the user. For example, an answer cannot consist of just some property that contains the keywords of the query, without showing the context.

The third drawback pertains to ranking, which must take into account textual relevance (as well as some other factors). In the naive approach, many nodes have only a small amount of text, making it hard to determine their relevance to a given query.

In comparison to the naive approach, the OCP model dictates *fat* nodes. That is, an object or an explicit connector is represented by a node that contains all of its properties. Consequently, we get the following advantages. First, a data graph is not unduly large, which improves efficiency. Second, relevance is easier to determine, because all the text pertaining to an object or an explicit connector is in the same node. Third, the GUI of [1] is effective, because it does not clutter the screen with too many nodes or unnecessary stuff. In particular, the default presentation of an answer is condensed and only shows: types and names of objects; types of explicit connectors; and properties that match some keywords of the query. The user can optionally choose an expanded view in order to see all the properties of the displayed nodes, when additional information about the answer is needed. Since all the properties are stored in the nodes that are already shown, this can be done without any delay. Furthermore, the GUI of [1] visualizes the

conceptual distinction between objects and connectors, which makes it much easier to quickly grasp the meaning of an answer.

4 Constructing Data Graphs

4.1 Relational Databases to Data Graphs

The naive approach for transforming a relational database into a data graph (e.g., [3]) does not distinguish between objects and connectors. For each tuple t , a node v_t is created, such that the relation name of t is the type of v_t . An edge from v_{t_1} to v_{t_2} is introduced if tuple t_1 has a foreign key that refers to t_2 . Finally, opposite edges are also added. In this section, we describe more elaborate rules that create a data graph with fat nodes and a clear distinction between objects and connectors.

As a matter of terminology, when we say “foreign key F ,” we mean that the foreign key consists of the set of attributes F . A foreign key F is transformed to a connector. Whether that connector is implicit or explicit depends on the names of the attributes comprising F . For example, suppose that there is a relation named `Student`. If the attribute `student` is a foreign key that points to that relation, then it can be transformed to an implicit connector. However, if the attribute `grader` points to the relation `Student`, it means that the foreign key corresponds to an entity that has a special role and is not just an ordinary student. In this case, the translation should create an explicit connector of type `grader`.

The above example serves as a motivation for the following definition. Suppose that F is a foreign key that refers to a relation R . Let P be the set consisting of the attributes of the primary key of R and R itself (i.e., the name of the relation). We say that the foreign key F is *insignificantly named* if $F \subseteq P$; otherwise, it is *significantly named*. For example, let $F = \{\text{student}\}$ be a foreign key that refers to the relation `Student` that has the primary key `id`; hence, $P = \{\text{Student}, \text{id}\}$. F is insignificantly named, because $F \subseteq P$. In practice, it is sufficient that F is similar (rather than strictly equal) to a subset of P . For example, if $F = \{\text{student_id}\}$, then we still deem F insignificantly named. (Due to a lack of space, we do not discuss how to test such similarity automatically.) If $F = \{\text{grader}\}$, then $F \not\subseteq P$ and hence F is significantly named.

Given a relation R , we transform its tuples to objects and connectors according to one of the following four cases.

1. The primary key of R does not include any foreign keys.
2. The primary key K of R includes a single foreign key F and either K has some attributes in addition to those of F or F is significantly named.
3. The primary key K of R is a combination of at least two foreign keys and possibly some other attributes.
4. The relation R has exactly one foreign key F , which is insignificantly named and is also the primary key.

In Case 1, 2 and 3, a tuple of R is an entity, a weak entity and a relationship, respectively. In these cases, we do the following. Each tuple t of R is transformed to a node v_t . In the first two cases, v_t is an object. In the third case, v_t is either an explicit

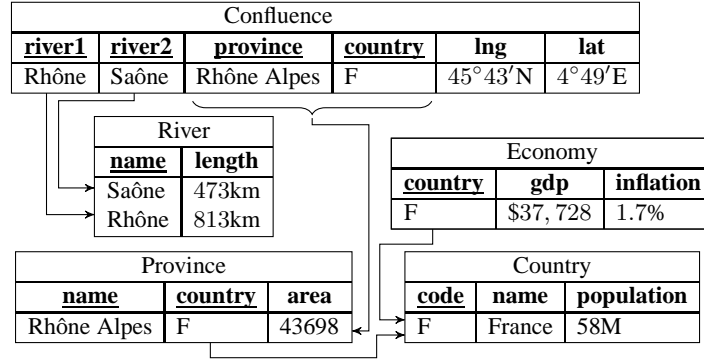


Fig. 6. A snippet of the Mondial RDB

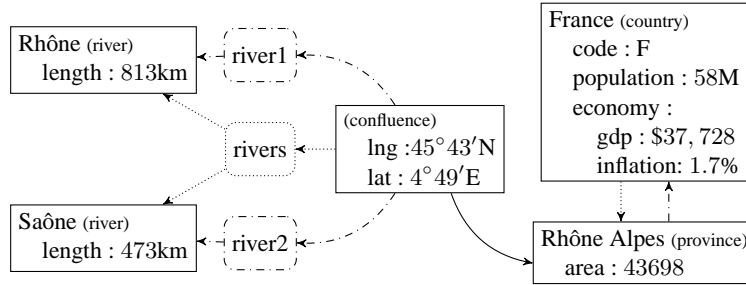


Fig. 7. A data graph constructed from the Mondial RDB

connector or an object, according to the following rule. If all the foreign keys of R are insignificantly named, then v_t is an explicit connector; otherwise, we make v_t an object (to avoid the creation of two explicit connectors that are adjacent).⁴

The type of v_t is R (i.e., the name of the relation). The properties of v_t are all the attributes of t that do not belong to foreign keys. If v_t is an object, its name is chosen to be the value of an appropriate property (e.g., title, name, etc.). In particular, we prefer a meaningful name over some meaningless id, even if the former does not uniquely identify the object.

In addition, for each foreign key F of t , we do the following. Let $v_{t[F]}$ be the object corresponding to the tuple referenced by $t[F]$ (i.e., the value of t for F). If F is insignificantly named, we add a directed edge e from v_t to $v_{t[F]}$ (note that e is an implicit connector if v_t is an object). Otherwise (i.e., F is significantly named), we create an explicit connector c_t^F of type F and add the directed edges (v_t, c_t^F) and $(c_t^F, v_{t[F]})$.

In Case 4, the relation R is an *auxiliary table* that provides additional information about the entities referenced by F . We transform a tuple t of R to a nested property of

⁴ Even when all the foreign keys of R are insignificantly named, v_t has to be an object if some other relation references tuples of R . However, since R represents a set of relationships (rather than entities), this possibility is unlikely to occur.

the object $v_{t[F]}$. The top-level property is R (i.e., the name of the relation) and it nests all the attributes of t that do not belong to F .

As an example, Figure 6 shows a snippet of the Mondial relational database. In each relation, the attributes of the primary key are underlined and arrows show foreign keys. We now explain how to construct the data graph of Figure 7. In this section, the dotted (implicit and explicit) connectors of Figure 7 should be ignored.

There are two relations, namely, `River` and `Country`, that satisfy the condition of Case 1. For each one of their tuples, Figure 7 has an object. Note that `France` is chosen to be the name of an object, although it is not the value of the primary key. The relation `Economy` satisfies the condition of Case 4. Therefore, its only tuple becomes the nested property `economy` of `France`.

The relation `Province` satisfies the condition of Case 2. Hence, the object `Rhône Alpes` of type `province` is in Figure 7; its only other property is `area`. The foreign key of `Province` is insignificantly named, so we add an implicit connector from `Rhône Alpes` to `France` that is shown as a dash-dotted arrow. Note that `country` is not a property of the object `Rhône Alpes`, because it belongs to a foreign key.

The relation `Confluence` of Figure 6 satisfies the condition of Case 3. Two out of the three foreign keys included in its primary key (i.e., `river1` and `river2`) are significantly named. Hence, the single tuple of `Confluence` is an object; however, there is a lack of an attribute that can serve as the name of that object. For each of the two significantly named foreign keys, we add an explicit connector, which is depicted using dash-dotted shapes (i.e., two arrows and a rectangle). The third foreign key comprises two attributes (`province` and `country`) and is insignificantly named. So, we add an implicit connector (shown as a solid arrow) from the object `confluence` to the object `Rhône Alpes`.

The above example shows that a constructed data graph could deviate from the original OCP model (of Section 2.1) in the following way. There is an object (of type `confluence`) without a name. It could be argued that this object should really be a connector. However, the result would be a data graph with adjacent connectors, which makes it harder to quickly grasp the meaning of answers having them. Moreover, a confluence actually corresponds to a real-world entity. In the RDB of Figure 6, it is a weak entity. So, we can create a name by concatenating the values of some primary-key attributes (e.g., `Rhône` and `Saône`).

The *original* edges are those created by the above transformation. We also add opposite edges (i.e., in the reverse direction), because the semantic of foreign keys is inherently undirected.

4.2 From XML to Data Graphs

An XML document is a rooted hierarchy of *elements*. Each element can have any number of *attributes*. Three special types of attributes are `ID`, `IDREF` and `IDREFS`. An attribute of the first type has a value that uniquely identifies its element. The last two types serve as references to other elements. For an attribute defined (in the DTD) as `IDREF`, the value is a single ID (of the referenced element); and if an attribute is defined as `IDREFS`, its value is a set of IDs. In our terminology, a *reference* attribute is

one defined as either IDREF or IDREFS. An attribute is *plain* if it is neither ID, IDREF nor IDREFS.

In XML lingo, an element has a *name* that appears in its tag (e.g., `<city>`). To avoid confusion, we call it the *type* of the element, because it corresponds to the notion of a type in the OCP model

In this section, we describe how to transform an XML document to a data graph. We assume that the document has a DTD and use it in the transformation. As we shall see, the DTD provides information that is essential to constructing the data graph. Conceivably, this information can also be gleaned from the document itself. However, if the document does not conform to a reasonable DTD, the resulting data graph (similarly to the document itself) is likely to be poorly designed. By only assuming that there is a DTD (as opposed to an XML schema), we make our transformation much more applicable to real-world XML documents.

Similarly to Section 4.1, we now define the concept of “significantly named;” we do it, however, for reference attributes, rather than foreign keys. Consider an attribute A that is defined as IDREF. A DTD does not impose any restriction on the type E of an element that can be referenced by the value of A . In a given XML document, A (i.e., its name) and E could be the same (e.g., `teacher`). If so, we say that A is an *insignificantly named* reference attribute. In the constructed data graph, the reference described by A can be represented by an implicit connector. If the opposite holds, namely, A and E are different, then we say that A is a *significantly named* reference attribute. In this case, the constructed data graph should retain A as the type of an explicit connector.

If attribute A is defined as IDREFS, then it is insignificantly named if all the IDs (in the value of A) are to elements of a type that has the same name as A ; otherwise, it is significantly named.

Whether a reference attribute is significantly named depends on the given XML document (and not just on the DTD). It may change after some future updates. As a general rule, we propose the following. It is safe to assume that a reference attribute A is significantly named if there is no element of the DTD, such that its type is the same as A . In any other case, it is best to get some human confirmation before deciding that a reference attribute is insignificantly named.

Let E_1 and E_2 be element types. We say that E_2 is a *child element type* of E_1 if the DTD has a rule for E_1 with E_2 on its right side. In this case, E_1 is a *parent element type* of E_2 .

Rudimentary rules for transforming an XML document to a data graph were given in [19]. However, they are applicable only to simple cases. Next, we describe a complete transformation that consists of two stages. We assume that prior to these two stages, both the DTD and the XML document are examined to determine for each reference attribute whether it is significantly named or not.

In the first stage, we analyze the DTD and classify element types as either objects, connectors or properties. This also induces a classification over the elements themselves. That is, when a type E is classified as an object, then so is every element of type E (and similarly when E is classified as a connector or a property). In the second stage, the classification is used to construct the data graph from the given XML document. The

first stage starts by classifying all the element types E that satisfy one of the following *base rules*.

1. If E does not have any child element type and all of its attributes are plain, then E is a property.
2. If E has an ID attribute or a significantly named reference attribute, then it is an object.
3. If E has neither any child element type nor an ID attribute, but it does have some reference attributes and all of them are insignificantly named, then E is a connector.

As an example, consider the DTD of Figure 8. Base Rule 2 implies that the element types `country`, `province`, `river` and `confluence` are objects, because the first three have an ID attribute and the fourth has a significantly named IDREFS attribute (i.e., `rivers`). No base rule applies to `economy`. By Base Rule 1, all the other element types are properties.

Next, we find all the element types that should be classified as properties by applying the following recursive rule. If (according to the DTD rules) element type E only has plain attributes and all of its child element types are already classified as properties, then so is E . It is easy to show that a repeated application of this recursive rule terminates with a unique result.

Continuing with the above example, a single application of the recursive rule shows that `economy` is a property, because all of its child elements have already been classified as such by Base Rule 1.

Now, we apply the following generalization of Base Rule 3. If E does not have an ID attribute, all of its child element types are classified as properties, and it has some reference attributes and all of them are insignificantly named, then E is a connector.

We end the first stage by classifying all the remaining element types as objects, and then the following observations hold. First, if an element type is classified as a property, then so are all of its descendants. Second, the classification (when combined with the construction of the data graph that is described below) ensures that a connector is always between two objects. Third, if an element type is classified as a connector, then it has some reference attributes and all of them are insignificantly named.

In the second stage, we transform the XML document to a data graph. At first, we handle PCDATA as follows. If an element e (of the document) includes PCDATA as well as either sub-elements or attributes, then we should create a new attribute having an appropriate name (e.g., `text`) and make the PCDATA its value. This is not needed if e has neither sub-elements nor attributes, because in this case, e becomes (in the data graph constructed below) a non-nested property, such that the element type of e is the name of that property and the PCDATA is its value.

Now we construct the data graph as follows. For each element e , such that e is not classified as a property, we generate a node n_e . This node is either an object or a connector (and hence an explicit one) according to the classification of e . The type of n_e is the same as that of e . If n_e is an object, we should choose one of its properties (which will be created by the rules below) as its name. As usual, we prefer a property (e.g., `title`) that describes the meaning of n_e , even if it is not a unique identifier. For each n_e , we create properties and add additional edges and nodes by applying the following six *construction rules*.

```

<ELEMENT country (name,population,
economy,province)>
<ATTLIST country (code ID #REQUIRED
area CDATA #IMPLIED)>
<ELEMENT economy (gdp,inflation)>
<ELEMENT province (name,area)>
<ATTLIST province (id ID #REQUIRED)>
<ELEMENT river (name,length)>
<ATTLIST river (id ID #REQUIRED)>
<ELEMENT confluence (lng,lat)>
<ATTLIST confluence
rivers IDREFS #REQUIRED)
province IDREF #REQUIRED)>
<ELEMENT name (#PCDATA)>
<ELEMENT population (#PCDATA)>
<ELEMENT gdp (#PCDATA)>
<ELEMENT inflation (#PCDATA)>
<ELEMENT area (#PCDATA)>
<ELEMENT length (#PCDATA)>
<ELEMENT lng (#PCDATA)>
<ELEMENT lat (#PCDATA)>

```

Fig. 8. DTD snippet of Mondial

```

<country code="F" area="547030">
  <name>France</name>
  <population>58M</population>
  <economy>
    <gdp>$37,728</gdp>
    <inflation>1.7%</inflation>
  </economy>
  <province id="prov-France-25">
    <name>Rhône Alpes</name>
    <area>43698</area>
  </province>
</country>
<river id="riv-Saone">
  <name>Saône</name>
  <length>473</length>
</river>
<river id="riv-Rhone">
  <name>Rhône</name>
  <length>813</length>
</river>
<confluence
  province="prov-France-25"
  rivers="riv-Saone riv-Rhone">
  <lng>45°43'N</lng>
  <lat>4°49'E</lat>
</confluence>

```

Fig. 9. XML snippet of Mondial

1. Every plain attribute of e is a property of n_e .
2. For each child p of e , such that p is classified as a property, the subtree (of the given document) that is rooted at p becomes a property of n_e . Note that this property is nested if p has either plain attributes or descendants of its own. Also observe that element types and attribute names appearing in p become names of properties nested in n_e .
3. For each child o of e , such that o is classified as an object (hence, so is e), we add an edge from n_e to n_o (which is the node created for o).
4. For each child c of e , such that c is classified as a connector, we add an edge from n_e to n_c . Observe that if such a c exists, then e is classified as an object and n_c is the node of the explicit connector corresponding to c .
5. For each reference attribute R of e , we create new connectors or add edges to existing ones, according to the following two cases. First, if R is insignificantly named, then for each object o that (the value of) R refers to, we add an edge from n_e to o . Note that this edge is an implicit connector if n_e is an object; otherwise, it is part of the explicit connector n_e .

The second case applies when R is significantly named. In this case, the classification rules imply that n_e is an object. We first create a node n_r , such that its only incoming edge is from n_e . This node represents an explicit connector that gets the name of attribute R as its type and has no properties. In addition, for each object o that (the value of) R refers to, we add an edge from n_r to o .

Figure 7 without the dash-dotted (but with the dotted and solid) parts shows the data graph created from the XML document of Figure 9 with the DTD of Figure 8. The two differences from the relational transformation of Section 4.1 are the following. First, the implicit connector between France and Rhône Alpes is from the former to the

latter (because the latter is a child of the former). Second, there is only one explicit connector `river` instead of `river1` and `river2`.

We divide the original edges (i.e., those created by the above transformation) into two kinds. The *hierarchical edges* are those created by Construction Rule 3. They are implicit connectors that reflect the parent-child relationship between XML elements. The *reference edges* are the ones introduced by Construction Rule 5 (i.e., due to reference attributes). Construction Rule 4 creates edges due to the element hierarchy, but they enter nodes of explicit connectors; hence, we also refer to them as reference edges.

As in the relational case (see Section 4.1), we add opposite edges. However, our experience indicates that even if it is done just for the reference edges (i.e., no opposite edges are added for the hierarchical ones), we generally do not miss meaningful answers to queries. Furthermore, as we show in the next section, a strategy that works well is to assign higher weights to opposite edges than to original ones. In this way, relevant answers are likely to be generated first without having too many duplicates early on.

5 A Comparison

In this section, we compare the data graphs produced from relational and XML data sources. At first, we describe the example about students, courses and lecturers that will be used.

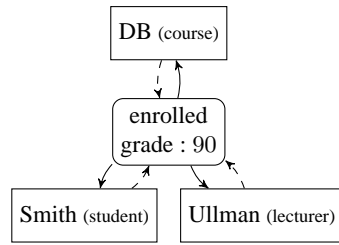


Fig. 10. Data graph from RDB with one ternary relationship

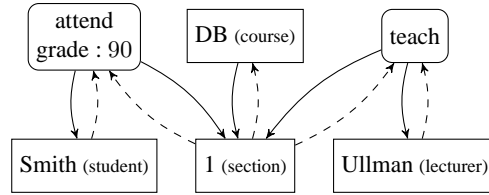


Fig. 11. Data graph from RDB with binary relationships

We abbreviate words by their first letter as follows: S(tudent), C(ourse), L(ecturer), E(nrolled) and G(rade). The three entity types student, course and lecturer have relations denoted by S , C and L , respectively. The attributes of those relations are not important. We only assume that each entity has a key and a name. By a slight abuse of notation, for each of these three relations, we will use its name also for denoting its key. Hence, the relationship between students, courses and lecturers is described by the relation $E(\underline{S}, \underline{C}, \underline{L}, G)$, where the attributes of the key are underlined. The data graph constructed according to Section 4.1 is given in Figure 10 (assuming that each relation has a single tuple). Opposite edges are shown as dashed arrows.

The relation $E(\underline{S}, \underline{C}, \underline{L}, G)$ involves three entity types, because a course may have several lecturers and not all of them teach every student attending the course. If courses

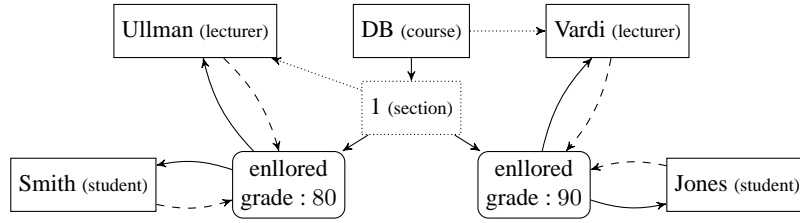


Fig. 12. Data graph from XML

are divided into sections, such that each one has its own lecturer(s), we can decompose E into two binary relationships. To incorporate sections, we will use the following abbreviations: A(ttend), T(each) and Sec(tion); that is, section is abbreviated by its first three letters. Now, the relation $E(\underline{S}, \underline{C}, \underline{L}, G)$ is replaced with $A(\underline{S}, \underline{C}, \underline{Sec}, G)$ and $T(\underline{C}, \underline{Sec}, \underline{L})$. Note that a section is a weak entity, and each of the new relations has two foreign keys, where one of them consists of two attributes (i.e., C and Sec) that together uniquely identify a section of a course; that is, the value of Sec is just a number, such as 1, 2, etc. The data graph for the five relations S , C , L , A and T is given in Figure 11.

The data graph produced from XML is shown in Figure 12 (the XML document and its DTD are not shown due to a lack of space). Figure 12 has more data (e.g., two lecturers) than Figures 10 and 11 to illustrate some points later. It has opposite edges (depicted as dashed arrows) only for reference (but not hierarchical) edges. One rectangle (for Section 1) and two edges are dotted. They are additions to the data graph that will be explained later. For now, they should be ignored (hence, there is an edge from DB directly to each enrolled connector).

To show the differences between the three data graphs, we consider the query $\{Student, Lecturer\}$. On the data graph of Figure 10, one answer has only original edges, an enrolled node as the root, and Smith and Ullman as the leaves. This answer is likely to be generated early on for two reasons. First, it is as small as can be (i.e., only three nodes). Second, it benefits from a strategy of assigning higher weights to opposite edges than original ones. In more detail, some algorithms (e.g., [3, 10, 14]) enumerate answers in an order that is correlated with increasing weight. If the most relevant answers are likely to have only original edges, then those algorithms would find them early on when opposite edges have higher weights.

For each course in which Smith is taught by Ullman, we would get an answer with an enrolled node as the root, and those two as the leaves. Either we remove duplicates by treating enrolled nodes according to their type, rather than id (i.e., all of them are identical to one another), or we should add the course so that users can grasp how seemingly duplicate answers are different from one another. In this case, we need to augment each answer with only one additional node, namely, the course object (pointed to by the enrolled connector). This requires adding only one node to each answer. The main drawback of the data graph of Figure 10 is a large number of enrolled connectors, since they represent a ternary relationship.

On the data graph of Figure 11, the answer with Smith and Ullman as the leaves must use a mixture of original and opposite edges, and has five nodes. We need to

add a sixth node if we want to show how duplicates are different from one another. In comparison with Figure 10, answers are larger implying that it takes longer to find them. Moreover, the strategy of assigning higher weights to opposite edges than original ones is not effective, because the most relevant answers have both types (so this approach would not help in generating them early on).

Next, we consider the data graph of Figure 12, which is obtained from XML. At first, we ignore the dotted rectangle and two edges. The answer that Ullman teaches Smith consists of only original edges and three nodes; a fourth one is needed to show the course. This is the same as in Figure 10, but the number of `enrolled` connectors is smaller (i.e., equal to the number of `attend` connectors in Figure 11, where `teach` connectors are also used). The main advantage of Figure 12, however, is its flexibility. If a lecturer teaches all the students enrolled in the course (which is likely to be true in many cases), then it is sufficient to have a connector from `course` to `lecturer`, such as the dotted edge from DB to Vardi (i.e., no need for edges between that lecturer and the `enrolled` connectors of students attending the course). Now, the subtree `Jones ← enrolled ← DB → Vardi` is the answer that Vardi teaches Jones. It consists of four nodes and already shows the course, which means that duplicates cannot occur. If we introduce sections (e.g., the dotted rectangle) and add edges to their lecturers (e.g., the dotted arrow pointing to Ullman), the subtree `Ullman ← 1 → enrolled → Smith` is the answer that Ullman teaches Smith. It consists of four nodes, and a fifth one should be added to show the course.

To summarize, a data graph obtained from XML has the following advantages over one constructed from a relational database.

1. Answers have an equal or smaller number of nodes when the same information (e.g., sections) is represented in both cases.
2. Relevant answers are more likely to use only original edges.
3. The data graph requires fewer nodes to represent ternary relationships (e.g., `enrolled`), because of the XML hierarchy.
4. The biggest advantage is heterogeneity:
 - It is sufficient to have sections only in courses that have more than one of them.
 - We can directly link a lecturer to a course, section or individual students depending on how she is assigned.
 - Thereby, we reduce the size of the data graph, give rise to fewer duplicates, and make answers more meaningful, because they show how the lecturer is assigned.

6 Conclusions

We showed that the OCP model is an effective conceptual basis for constructing data graphs. Using it, we developed transformations for generating data graphs from RDB and XML. These transformations are quite elaborate and provide much better results than the ad hoc methods that have been used in the literature thus far. In particular, the produced data graphs are better in terms of both efficiency (i.e., answers are generated more quickly) and effectiveness (i.e., the most relevant answers are produced early on).

It should be emphasized that the presented transformations are based on the principle of creating fat nodes (as explained in Section 3) and avoiding redundancies (e.g., due to insignificantly named references). Thus, they are applicable and useful (in most if not) all cases, regardless of how answers are generated or ranked.

We showed that XML is the preferred starting point for constructing data graphs. However, we need to better understand how to create XML documents that yield the best possible data graphs. Toward this end, we plan to develop appropriate design rules for XML documents.

Due to space limitations, we did not discuss how to generate data graphs from RDF. In our experience, it is harder to do that than when starting with RDB or XML. An important principle of RDF is unique representation by means of URIs (uniform resource identifiers). As a result, RDF triples are highly fragmented (e.g., there could be a separate triple for storing each person's title, such as Dr., Mrs., etc.), which makes it hard to create a coherent data graph with fat nodes.

An interesting topic for future work is to how to construct data graphs from XML documents without DTDs.

References

1. Achiezra, H., Golenberg, K., Kimelfeld, B., Sagiv, Y.: Exploratory keyword search on data graphs. In: SIGMOD Conference (2010)
2. Bao, Z., Ling, T.W., Chen, B., Lu, J.: Effective XML keyword search with relevance oriented ranking. In: ICDE (2009)
3. Bhalotia, G., Hulgeri, A., Nakhe, C., Chakrabarti, S., Sudarshan, S.: Keyword searching and browsing in databases using BANKS. In: ICDE (2002)
4. Coffman, J., Weaver, A.C.: An empirical performance evaluation of relational keyword search techniques. *IEEE Trans. Knowl. Data Eng.* 26(1), 30–42 (2014)
5. Dalvi, B.B., Kshirsagar, M., Sudarshan, S.: Keyword search on external memory data graphs. *PVLDB* (2008)
6. Ding, B., Yu, J.X., Wang, S., Qin, L., Zhang, X., Lin, X.: Finding top-k min-cost connected trees in databases. In: ICDE (2007)
7. Golenberg, K., Kimelfeld, B., Sagiv, Y.: Keyword proximity search in complex data graphs. In: SIGMOD Conference (2008)
8. Golenberg, K., Kimelfeld, B., Sagiv, Y.: Optimizing and parallelizing ranked enumeration. *PVLDB* (2011)
9. Golenberg, K., Sagiv, Y.: Constructing data graphs for keyword search. In: DEXA (2016), presentation url: https://drive.google.com/open?id=0BxX7DI4NO_-vTFZaRGgzU25WdjQ
10. Golenberg, K., Sagiv, Y.: A practically efficient algorithm for generating answers to keyword search over data graphs. In: ICDT (2016)
11. Guo, L., Shao, F., Botev, C., Shanmugasundaram, J.: XRANK: Ranked keyword search over XML documents. In: SIGMOD Conference (2003)
12. He, H., Wang, H., Yang, J., Yu, P.S.: BLINKS: ranked keyword searches on graphs. In: SIGMOD Conference (2007)
13. Hristidis, V., Papakonstantinou, Y., Balmin, A.: Keyword proximity search on XML graphs. In: ICDE (2003)
14. Kacholia, V., Pandit, S., Chakrabarti, S., Sudarshan, S., Desai, R., Karambelkar, H.: Bidirectional expansion for keyword search on graph databases. In: VLDB (2005)

15. Kasneci, G., Ramanath, M., Sozio, M., Suchanek, F.M., Weikum, G.: STAR: Steiner-tree approximation in relationship graphs. In: ICDE (2009)
16. Li, G., Ooi, B.C., Feng, J., Wang, J., Zhou, L.: EASE: an effective 3-in-1 keyword search method for unstructured, semi-structured and structured data. In: SIGMOD Conference (2008)
17. Mass, Y., Sagiv, Y.: Virtual documents and answer priors in keyword search over data graphs. In: Proceedings of the Workshops of the EDBT/ICDT 2016 Joint Conference (2016)
18. Park, C., Lim, S.: Efficient processing of keyword queries over graph databases for finding effective answers. *Inf. Process. Manage.* 51(1), 42–57 (2015)
19. Sagiv, Y.: A personal perspective on keyword search over data graphs. In: ICDT. pp. 21–32 (2013)
20. Xu, Y., Papakonstantinou, Y.: Efficient LCA based keyword search in XML data. In: EDBT (2008)