

Converting XML Data To UML Diagrams For Conceptual Data Integration

Mikael R. Jensen Thomas H. Møller Torben Bach Pedersen
Department of Computer Science, Aalborg University
{mrj, thm, tbp}@cs.auc.dk

Abstract

The demand for data integration is rapidly becoming larger as more and more information sources appear in modern enterprises. In many situations a logical (rather than physical) integration of data is preferable since some data is inherently not suited for storing in a physically integrated data warehouse. Previous web-based data integration efforts have focused almost exclusively on the logical level of data models, creating a need for techniques focused on the conceptual level. Extensible Markup Language (XML) is fast becoming the new standard for data representation and exchange on the World Wide Web, e.g., in B2B e-commerce, making it necessary for data analysis tools to handle XML data as well as traditional data formats..

This paper presents algorithms for automatically constructing UML diagrams from XML data and shows how to use the diagrams for the conceptual design of (virtual) Data Warehouses based on web data. Unlike earlier work, the focus of this paper is the support of On-Line Analytical Processing (OLAP) systems based on web data.

1 Introduction

Integration of distributed data sources is becoming increasingly important as more business relevant data appear on the web, e.g., on B2B marketplaces, and enterprises cooperate more tightly with their partners, creating a need for integrating the information from several enterprises. The *data warehousing* approach dictates a physical integration of data, enabling fast evaluation of complex queries, but also demanding great effort in keeping the data warehouse up to date. However, some kinds of data are impossible, or not attractive, to physically store in a data warehouse, either because of legislation or because of the nature of the data, e.g., rapidly dynamically evolving information. Enabling integrated use of such data requires a logical rather than physical integration.

XML is a meta language used to describe the structure and content of documents. XML is increasingly used for data exchange on the Web. The application of XML as a standard exchange format for data publicly available on the Web makes it attractive to use for building data warehouses from web data. Previous approaches for integrating web-based data, particularly in XML format, have focused almost exclusively on data integration at the *logical* level of the data model, creating a need for techniques that are usable at the conceptual level which is more suitable for use by system designers and end users. The most wide-spread conceptual model is the Unified Modeling Language (UML) [13].

The contribution of this paper is twofold. First, algorithms for automatically constructing UML diagrams from XML data are presented, enabling fast and easy graphical browsing of XML data sources on

the web. The algorithms capture important semantic properties of the XML data such as precise cardinalities and aggregation (containment) relationships between the data elements. These precise properties would be lost if an intermediate translation (XML to relational, followed by relational to UML) was used. Second, an architecture for integrating XML data at the conceptual level is presented. The architecture also supports relational data sources and is thus well suited to build data warehouses which are based partly on in-house relational data and partly on XML data available on the web. The focus of the data integration effort is the support of OLAP systems, meaning that some design choices are different than in a more general setting, e.g., the order of data in the XML documents is considered unimportant as we are interested in *aggregation* queries that destroy document order anyway. We assume that the XML data has a Document Type Description (DTD) to specify its structure. More powerful frameworks for describing XML documents, most notably XML Schema [19] exist, but none are yet official W3C recommendations. Also, only DTDs are currently widely used in real-world XML documents. Also, even if XML Schema is used to describe the XML data, the process of designing the integrated database will still benefit significantly from the graphical notation found in UML (and not in XML Schema).

We believe this paper to be the first to present algorithms that automatically construct UML diagrams directly from XML data, retaining important semantic information. Also, we believe it to be the first to consider the important issue of integration of web-based data sources at the conceptual level and to focus on OLAP systems.

A number of systems allows integration of structured and/or semi-structured data on the web. The systems YAT [1], Strudel [4], TSIMMIS [6], and Ozone [10] all use a semi-structured/XML data model, Garlic [16] uses an object data model, Cohera [7] uses the relational (SQL) data model, while Clio [12] supports combined object-relational and XML data. These systems all focus on web data integration at the *logical* level. In comparison, this paper focuses on integration of web data on the conceptual level, supporting the *design* process better. Also, our focus is on OLAP systems rather than generic data integration. However, we do not consider data integration at the logical level to be unnecessary, as such systems are used for the *implementation* of the integrated databases. Indeed, a DW conceptually designed using our approach could be implemented using these systems. The issue of converting XML data into relational data is described in [3, 5, 15]. In comparison, this paper aims at capturing more semantics of the XML data by using a higher-level conceptual model rather than providing efficient querying of XML data using RDBMSes. Database design tools such as ERwin [2] can “reverse-engineer” relational structures into UML diagrams. However, semantic information such as precise cardinalities and containment type relationships invariably get lost in the translation process. In comparison, this paper translates XML data directly into UML diagrams, keeping this valuable information.

The remainder of the paper is organized as follows. Section 2 defines algorithms for automatically generating a UML diagram from an XML DTD. Section 3 presents an integration architecture that enables XML and relational data at the conceptual level for constructing a “virtual” DW. Section 4 summarizes and points to topics for future research.

2 Generating UML Diagrams From XML Data

This section describes algorithms for automatically generating UML diagrams from XML data, enabling fast and easy browsing of XML data available on the Web. The structure of XML data is visualized by a UML diagram derived from a DTD describing the XML data source. Furthermore, suggestions for

determining types for XML elements are briefly presented. The purpose of the algorithms is to build a UML diagram which is at least as general as the document described by a DTD, meaning that XML data conforming to the DTD must be expressible in terms of the derived UML diagram. The algorithms generate UML according to the OMG Unified Modeling Language Specification version 1.3 [13]. Only a subset of the modeling constructs described in this specification are used, namely classes, typed attributes, and cardinality- and role-specified associations and aggregations. The algorithms are designed such that the structure of the original XML document is preserved and visible in the UML diagram.

The UML diagram has been chosen as the way of modeling and visualizing the DTD because UML is a standardized conceptual data modeling language which is familiar to many designers and end users, and because UML is powerful enough to express any document described by a DTD. The diagram is a conceptual easy understandable way of visualizing the structure of the XML document and allows the user quickly to get a feeling of the content and the structure of the data available. This supports easy browsing of data made available through XML on the web. The algorithms for deriving a UML diagram from a DTD can in addition to the nesting constructs of XML also handle recursion between elements (both direct and indirect recursion by either nesting or reference) as well as references between different elements (ID/IDREF(S) constructs).

It is assumed that the XML data is described by a DTD and that it is valid according to this DTD. We also assume that all ID-references (IDREFS referring to some IDs) from a given element is a reference to elements of one type only. This means that for an element having multiple references defined (using an IDREFS attribute), each reference must be to the same element type. This assumption has been made since it in common data models, like the UML model, E-R model, and relational model, is generally impossible to model references to multiple unspecified elements¹. Also, this only occurs very rarely in well-designed XML documents. It is furthermore assumed that it is possible to determine the relationship between any two elements having an ID/IDREF(S) relationship. This means that for every element having an IDREF or IDREFS attribute defined, it is possible to determine the element type to which a reference exist. The XML 1.0 Recommendation [18] for XML data is used throughout this paper

2.1 Pre-processing of the DTD

Before generating the UML diagram corresponding to the document described by the DTD, the DTD is simplified in order to make further processing easier. This simplifying process is here termed DTD pre-processing. It is not the intention to create a DTD that is equivalent to the original DTD. Instead the transformations creates a DTD which can be “less strict” than the original DTD, meaning that any XML document described by the original DTD is “at least” described by the transformed DTD. This assures that any document conforming to the original DTD also conforms to the transformed DTD.

All that matters in an XML document is the position of an element relative to its siblings and the parent-child relationship between elements, and it is therefore safe to apply the transformations since they preserve these conditions. The pre-processing makes the following changes to the DTD:

- Removes elements not accessible from the document root
- Removes elements declared EMPTY which have not declared an ATTLIST
- Simplifies element content specifications

¹The only way is to make zero-to-many relationships between all entities in the model, a rather unsatisfactory solution.

The first step towards simplifying the DTD is to remove all elements not accessible from the root. It is safe to apply this procedure since any XML document conforming to the DTD cannot make use of this type of elements. [18] All elements declared EMPTY which does not have an ATTLIST declared can safely be removed from the DTD. It is possible to make use of this type of elements in an XML document, but the element can contain no data and are therefore of no interest. [18]

Most of the complexity of a DTD stems from the complex definition of element content specifications. By applying certain transformations to a DTD it is possible to simplify the DTD element content specifications. The transformations used in the pre-processing are a slight change of the transformations presented in [15], consisting in that “+” operators are not transformed into “*” operators. This is not done since it is possible to capture the “+” operator semantics in the UML model. Only a subset of the transformations are shown here.

Flattening transformations	Simplification transformations	Grouping transformations
$(e_1, e_2)^* \rightarrow e_1^*, e_2^*$	$e_1^{**} \rightarrow e_1^*$	$\dots, a^*, \dots, a^*, \dots \rightarrow a^*, \dots$
$(e_1, e_2)? \rightarrow e_1?, e_2?$	$e_1^{*?} \rightarrow e_1^*$	$\dots, a^*, \dots, a?, \dots \rightarrow a^*, \dots$
$(e_1 e_2) \rightarrow e_1?, e_2?$	$e_1^{?*} \rightarrow e_1^*$	$\dots, a?, \dots, a^*, \dots \rightarrow a^*, \dots$
	$e_1{??} \rightarrow e_1?$	$\dots, a?, \dots, a?, \dots \rightarrow a^*, \dots$
		$\dots, a, \dots, a, \dots \rightarrow a^+, \dots$

The first set of transformations convert a nested definition into a flat representation, meaning that the binary operators “,” and “|” do not appear inside any operator. The second set of transformations reduces many unary operators to a single unary operator. The last set of transformations group sub-elements having the same name. Note that some of the transformations destroy document order. This is unproblematic since we are focusing on OLAP queries that do not query the document order of the detail data but rather aggregate several detail data items into one aggregate result, thus destroying document order anyway.

Example 2.1 Simplification of Element Content Specification

Consider the specification $\langle !ELEMENT\ x\ ((b|c)^*, a, (d?(e?, (b, b)^*))^*, a) \rangle$. By applying the transformations it can be reduced to $\langle !ELEMENT\ x\ (b^*, c^*, a^+, d^*, e^*) \rangle$.

2.2 DTD and UML Data Model

This section defines the formal data model specifications used to describe the DTD and the UML diagram components. Both the DTD and the UML data model are defined by means of an environment based on sets. The two data model specifications are defined in order to give a precise description of the generation and modification of a UML diagram.

Notation In order to ease the readability of the algorithms presented, the following notation will be used. All sets are written in *italics* with initial capital letter. A capital letter is used to separate words of a set, e.g. *MySet*. All elements of sets are written in *italics* with small letters. An underscore character is used to separate words of an element, e.g. *my_element*. All functions are written in *sans serif* with small letters. A capital letter is used to separate words of a function, e.g. *myFunction*. When a set consists of elements that are tuples square brackets (“[“ and “]”) are used to address a certain tuple

element, e.g. when an element is defined by $my_element = (my_tuple_element1, MyTupleElement2)$ then $my_element[my_tuple_element1]$ denotes the first element in the tuple. Furthermore, when a tuple consists of elements which are sets (like the second tuple from the former example) then a subscript will denote a certain element within this set, e.g. $my_element[MyTupleElement2_j]$ denotes element j in the set $MyTupleElement2$. We now define the sets needed to describe an XML DTD.

$$ElementNames = \{ e \mid e \text{ is an element name defined in the DTD} \} \cup \{ CDATA, PCDATA \}$$

$$AttributeNames = \{ a \mid a \text{ is an attribute name defined in the DTD} \}$$

$$DTDCContent = \{ (element_name, ElementContentSpec, AttContentSpec) \mid \\ element_name \in ElementNames \}$$

$$ElementContentSpec = \{ (element_name, modifier) \mid element_name \in ElementNames \wedge \\ modifier \in Modifiers \setminus \{ REQUIRED, IMPLIED \} \}$$

$$AttContentSpec = \{ (att_name, att_type, modifier) \mid att_name \in AttributeNames \wedge \\ att_type \in AttributeTypes \wedge modifier \in Modifiers \setminus \{ *, +, ?, 1 \} \}$$

$$Modifiers = \{ *, +, ?, 1, REQUIRED, IMPLIED \}$$

$$AttributeTypes = \{ ID, IDREF, IDREFS, CDATA \}$$

The set *ElementNames* contains all names of elements declared with *!ELEMENT* in the DTD including the pre-defined character data elements CDATA and PCDATA. *AttributeNames* contains all names of attributes declared with *!ATTLIST* in the DTD. The *DTDCContent* set is a set of three-tuples describing the contents and structure of the DTD. The first element in the tuple is the name of an element from the DTD. The second and third elements in the tuple are both sets which are defined below. The *ElementContentSpec* set describes the content specification for all elements defined in the DTD. The set contains a two-tuple for each element listed in the element content specification for some element in the DTD. A modifier is associated with each element in accordance to the modifier specified in the DTD. For each element declared in the DTD a list of attributes can be associated with the element. The set *AttContentSpec* contains three-tuples where the first element in the tuple is the name of the declared attribute, the second element is the type of the attribute and the third element is the modifier associated with the attribute. *Modifiers* is the set of modifiers used to describe the cardinality of elements and attributes specified in the DTD. All attributes defined in a DTD are typed. *AttributeTypes* is the set of attribute types.

Definition A *leaf element* is declared in the DTD to have character data, e.g., PCDATA, as content.

We now define the sets for describing the UML diagram. The set *UMLClasses* describes a complete UML diagram. It contains three-tuples where the first element in the tuple is the name of the UML class. The two other elements are both sets, where *AttributeList* is the list of attributes associated with this particular class and *PointsTo* is a list of classes that this class links to. The element r is a Uniform Resource Identifier (URI) uniquely identifying the data source.

$$UMLClasses = \{ (class_name, AttributeList, PointsTo) \mid class_name \in ElementNames \} \cup \\ \{ r \mid r \text{ is a Uniform Resource Identifier} \}, \text{ where}$$

$$AttributeList = \{ (att_name, modifier, data_type) \mid \\ modifier \in Modifiers \setminus \{ *, +, REQUIRED, IMPLIED \} \wedge \\ data_type \in DataTypes \}, \text{ and}$$

$$\begin{aligned}
PointsTo = & \{ (class_name, link_type, source_card, target_card, link_role) \mid \\
& class_name \in ElementNames \wedge link_type \in LinkTypes \wedge \\
& source_card, target_card \in Cardinalities \wedge \\
& link_role \in AttributeNames \cup \{NULL\}
\end{aligned}$$

The *AttributeList* set describes the name, modifier and data type of each attribute in a UML class. The modifier can either be 1 or ?, describing whether or not null values are allowed for this attribute. The *data_type* element holds type information for this variable. The *PointsTo* describes a relationship from one class to another class. *class_name* is the name of the class to which a link exist, *link_type* is the type of link and *source_card* and *target_card* is the cardinality specified for the source and the target of the link, respectively. *link_role* is the role of the link. A link role is a description of a link, containing a name and a direction. In this case only the name is used since the direction is implicit (the link is always from the class holding the *PointsTo* element to the class named in *class_name*). Only association type links are given a role, since the meaning of an aggregation type link is always clear (aggregation means “contained in”), and links of this type are thus assigned NULL as *link_role*.

$$\begin{aligned}
DataTypes &= \{ NUMERIC, DATE, TEXT, NULL \} \\
LinkTypes &= \{ AGGREGATION, ASSOCIATION \} \\
Cardinalities &= \{ 0..*, 1..*, 0..1, 1 \}
\end{aligned}$$

The *DataTypes* set describes four basic data types for attributes in the UML diagram. Data types are not applicable when generating the UML model, since the DTD contains no type information (see Section 2.6). The set *LinkTypes* describes the different types of relations that can exist between any pair of classes in the UML diagram. *Cardinalities* is the set of cardinalities used to describe the quantitative relationship between elements in the UML data model. The cardinalities of a relationship are given by specifying minimum and maximum cardinalities. The cardinality “1” is used as a shorthand for “1..1”.

Functions We now define the functions used in the conversion from DTDs to UML diagrams. All functions are defined by specifying the domain of definition and range. A “ \rightarrow ” denotes total functions whereas a “ \hookrightarrow ” denotes a partial function.

$$\begin{aligned}
target &: AttContentSpec \hookrightarrow ElementNames \\
cardinality &: Modifiers \cup AttContentSpec \rightarrow Cardinalities \\
parent &: UMLClasses \hookrightarrow UMLClasses
\end{aligned}$$

The *target* function is given an attribute specification and returns the name of an element. The function is used to pair IDREF(S) and IDs, and is only defined for the subset of *AttContentSpec* having an *att_type* of either IDREF or IDREFS. The *cardinality* function is given either an element from *Modifiers* or an element from *AttContentSpec* and returns the cardinality corresponding to the modifier or the cardinality of the attribute. An attribute declared as IMPLIED has cardinality 0..1, an REQUIRED attribute has cardinality 1, an IMPLIED IDREFS attribute has cardinality 0..*, an REQUIRED IDREFS attribute has cardinality 1..*. The modifiers *, + and ? has cardinality 0..*, 1..* and 0..1, respectively. The *parent* function is given an element from *UMLClasses* and returns the parent element of this element as specified

by the nesting relationship in the DTD. The function is only defined for the elements having exactly one parent.

2.3 DTD to UML Conversion Algorithm

The algorithm for generating the UML diagram has been divided into two parts. The first part, `generateUMLClass`, is a subfunction that generates one element in *UMLClasses* which corresponds to a complete class in the final UML diagram. The other part of the algorithm, `generateUML`, calls the `generateUMLClass` subfunction repeatedly until all the elements of *UMLClasses* have been generated. The algorithms generate only valid UML diagrams. The only UML components generated by the algorithm `generateUMLClass` are classes with zero or more attributes, cardinality-specified aggregations between classes and, role- and cardinality-specified associations between classes (see above for a definition of association roles). Aggregations are always constructed in direction from parent to child and associations are always constructed in direction from the referring element to the element bearing the corresponding ID.

In general the structure of the UML diagram generated by the algorithm corresponds to the tree-structure of the DTD where a UML class is generated for each element. All elements having a parent/child relationship in the DTD are connected by an aggregation relationship in the UML diagram and ID-references between elements are modelled as associations between the corresponding UML classes.

GenerateUMLClass Algorithm

- (1) `generateUMLClass($x_i \in DTDCContent$):`
- (2) Generate new element $e \in UMLClasses$ where
- (3) $e[class_name] = x_i[element_name] \quad \wedge$
- (4) $e[AttributeList] = \{(a, b, c) \mid a = x_i[AttContentSpec_j][att_name],$
- (5) $b = cardinality(x_i[AttContentSpec_j][modifier]), c = NULL \wedge$
- (6) $x_i[AttContentSpec_j][att_type] \notin \{IDREF, IDREFS\}\} \cup$
- (7) $\{(a, b, c) \mid a = \text{“value”}, b = 1, c = NULL \wedge x_i \text{ is a leaf element}\} \quad \wedge$
- (8) $e[PointsTo] = \{(d, e, f, g, h) \mid d = x_i[ElementContentSpec_j][element_name],$
- (9) $e = AGGREGATION,$
- (10) $f = 1,$
- (11) $g = cardinality(x_i[ElementContentSpec_j][modifier]),$
- (12) $h = NULL \quad \wedge$
- (13) $x_i \text{ is not a leaf element}\} \cup$
- (14) $\{(d, e, f, g, h) \mid d = target(x_i[AttContentSpec_j]),$
- (15) $e = ASSOCIATION,$
- (16) $f = cardinality(x_i[AttContentSpec_j]),$
- (17) $g = 0..*,$
- (18) $h = x_i[AttContentSpec_j][att_name] \quad \wedge$
- (19) $x_i[AttContentSpec_j][att_type] \in \{IDREF, IDREFS\}\}$

The subfunction `generateUMLClass` works as follows: The function is given an element $x_i \in DTDCContent$ and generates an element $e \in UMLClasses$ which is the representation of a UML class. $e[class_name]$ is set to the name of the element in the DTD to which x_i corresponds. $e[AttributeList]$ is a

set of attribute names, modifiers and data types for all the non IDREF and IDREFS type attributes defined for the UML class corresponding to x_i . The data type is not applicable since it cannot be determined at this point. If the element x_i is a leaf element an attribute of name “value” is added to the class to hold the element’s data. This means that all elements having their content defined as CDATA or PCDATA will have the attribute “value” to hold this character data.

$e[PointsTo]$ is a set of links to each of the elements within the content specification for the element corresponding to x_i . The link type and cardinality is determined by the functions `linkType` and `cardinality`. The cardinality *source_card* for links caused by attributes of types different from IDREF and IDREFS always defined as 1 since elements in XML are always uniquely nested. The cardinality *target_card* is for attributes of type IDREF and IDREFS always set to 0..* since it in XML documents described by DTDs not can be controlled how many elements are allowed to refer to an element having an ID attribute. A link can be from one element to itself indicating a recursive definition.

The link types are either aggregations or associations. Aggregations are used whenever a link to a nested element is established (modeling the “consists of” parent-child relationship) and associations are used when modeling an ID-reference relationship between two elements. A link role is supplied for each association type link, containing the name of the IDREF or IDREFS attribute causing the link. This helps resolving the conceptual meaning of association type links in the UML diagram, provided, of course, that the IDREF and IDREFS attributes are given meaningful names in the DTD.

GenerateUML Algorithm

- (1) `generateUML`:
- (2) $UMLClasses = \bigcup_i \{ \text{generateUMLClass}(x_i \in DTDContent) \}$
- (3) $UMLClasses = UMLClasses \cup \{r\}$

The algorithm works as follows: The set *UMLClasses* is generated by uniting the result of calling the `generateUMLClass` subfunction on every $x_i \in DTDContent$. *r* is the URI describing the location of the XML document.

2.4 UML Post-processing

To reduce the number of classes in the UML diagram *UMLClasses* is processed by the algorithm `post-ProcessUML`, seen below. The number of classes is reduced by removing all leaf elements having only one parent and no association relationships. The attributes and data contained within the leaf element is moved from the leaf element to its immediate parent, thereby making it possible to remove the now empty leaf element.

The algorithm works as follows: An element $e_i \in UMLClasses$ is a candidate for removal if e_i is a leaf element having one parent and no incoming or outgoing references. When a candidate element has been found, it is checked if the link from e_i ’s parent to e_i has a target cardinality of 1 or 0..1. If it has, e_i can be safely removed. The reason for only removing classes having this cardinality is that classes that are part of a “many” relationship are always modelled as separate classes in UML.

When an element e_i destined for removal has been located, the names of the attributes defined for the element are changed. The new attribute names are a concatenation of the name of the class, a dot (“.”) and the old attribute name (in the algorithm “ \circ ” is the concatenation operator). Finally, all the attributes

from the about-to-be removed class are given a modifier corresponding to the class' cardinality and are then copied to the parent.

Post-processing Algorithm

- (1) postProcessUML:
- (2) $\forall e_i \in UMLClasses$ where e_i is a leaf element having one parent and no references:
- (3) **if** $(parent(e_i)[PointsTo_j][class_name] == e_i[class_name]) \wedge$
- (4) $(parent(e_i)[PointsTo_j][target_card] \in \{ 1, 0..1 \})$ **then**
- (5) $\forall a_m \in e_i[AttributeList]$:
- (6) $a_m[att_name] = e_i[class_name] \circ "." \circ a_m[att_name]$
- (7) **if** $(parent(e_i)[PointsTo_j][target_card] == 0..1)$ **then**
- (8) $\forall z_k \in e_i[AttributeList]$: $z_k[modifier] = ?$
- (9) **endif**
- (10) $parent(e_i)[AttributeList] = parent(e_i)[AttributeList] \cup e_i[AttributeList]$
- (11) $parent(e_i)[PointsTo] = parent(e_i)[PointsTo] \setminus \{ parent(e_i)[PointsTo_j] \}$
- (12) $UMLClasses = UMLClasses \setminus \{ e_i \}$
- (13) **endif**

2.5 Example: Generating a UML Diagram from an XML DTD

This section illustrates through an example how the previous algorithms work. It is assumed that the document used in the example is located at <http://www.componentheaven.com/parts.xml>.

Consider the following DTD, based on the Electronic Component Information Exchange (ECIX) QuickData Architecture [17], a project dedicated to the design of standards for B2B technical information exchange of component information for traditional electronic components. The root element of this DTD is the `class` element.

```
<!ELEMENT class ((ec, device)*) >
<!ATTLIST class name CDATA #REQUIRED >
<!ELEMENT ec (unitprice, pincount, gatecount, textdesc?) >
<!ATTLIST ec id ID #REQUIRED usedWithin IDREFS #IMPLIED
name CDATA #REQUIRED >
<!ELEMENT device (textdesc, unitprice, device?) >
<!ATTLIST device id ID #REQUIRED name CDATA #REQUIRED >
<!ELEMENT unitprice (number, price) >
<!ELEMENT pincount (#PCDATA) >
<!ELEMENT gatecount (#PCDATA) >
<!ELEMENT textdesc (#PCDATA) >
<!ELEMENT number (#PCDATA) >
<!ELEMENT price (#PCDATA) >
```

When pre-processing the DTD according to the transformations described in Section 2.1 the element content specification for element “class” is reduced to `<!ELEMENT class (ec*, device*) >`.

We now proceed to construct the sets describing the DTD as defined above. Due to space constraints, only an overview is given here.

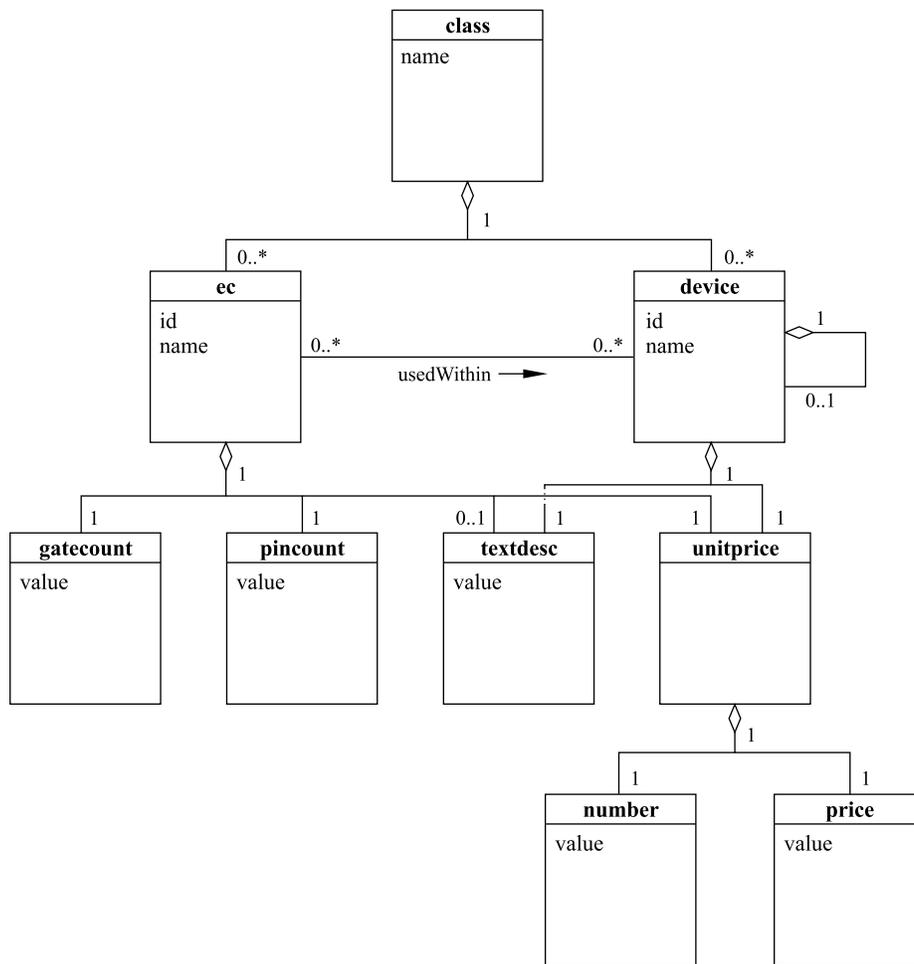


Figure 1: UML Diagram Produced by generateUML Algorithm

A UML diagram describing the DTD can be built directly from the sets describing the UML diagram. The UML diagram for the DTD in this example is shown in Figure 1. Notice the labelled arrow on the association type link from “ec” to “device”, which is a representation of the *link_role* element associated with this link. The labelled arrow indicates the direction and name of this link, making it easier to understand the relationship between the two classes.

The UML diagram in Figure 1 has one class for each element specified in the DTD. To reduce the total number of classes the `postProcessUML` algorithm is applied to the *UMLClasses* set. The UML diagram resulting from the application of `postProcessUML` is shown in Figure 2. As can be seen from the figure, four classes have been removed by the post processing algorithm. This is caused by the “in-lining” of simple leaf classes, resulting in a diagram that is more meaningful and easier to understand for humans.

2.6 Determining Attribute Data Types

In the generation of the UML diagram a set of attributes is associated with each of the classes as described above. Since the UML diagram is generated directly from the XML DTD it is not possible to determine data types for the attributes at that point, and they are therefore typed with `NULL`. However, before constructing a DW, data types has to be determined for the attributes related to the post-processed UML diagram. It is not the intention of this section to provide a final solution of how type information is

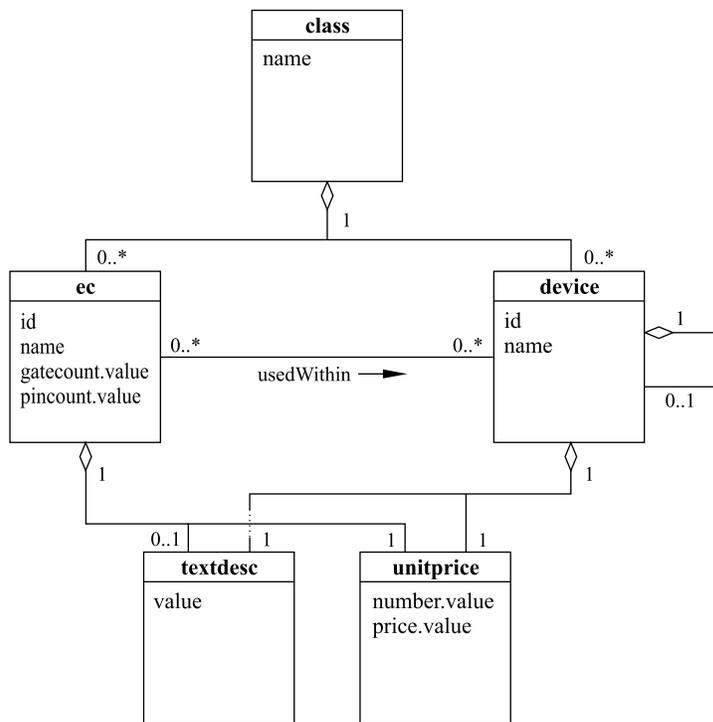


Figure 2: UML diagram after post-processing

assigned to attributes. Instead, some of the problems associated with type determination is presented along with brief suggestions on how to resolve these difficulties.

One of the problems in determining type information for the attributes, is that XML 1.0 is a typeless language, so no data type declaration exists in the DTD or in the XML document. [18] The type space for XML is simply string values [18], but this information is too vague to be of any practical use in a DW system (e.g. when wanting to use aggregation functions other than count). Therefore the type space is extended as described above by the set *DataTypes* (having members NUMERIC, DATE, etc.).

One way to assign type information to the attributes is to sample from the XML document. For each of the attributes in the UML diagram a number of samples can be done from the data contained within the corresponding XML tags, thereby making it possible with a certain probability to make a reasonable type determination.

Example 2.2 Sampling

As an example of this approach consider the XML fragment `<price>123</price>`. By using sampling, the attribute price could be assigned the type NUMERIC.

The designer can also assist in determining types. It is fair to assume that the designer has some a priori knowledge about the data used for the analysis. By presenting the attributes of concern, the designer can make type assignments to these attributes based on a priori knowledge. Another way to resolve types is to simply provide a lexical analysis of the element declarations in the DTD, in order to make qualified guesses of types based on certain keywords.

Example 2.3 Lexical Analysis

By using the approach based on lexical analysis an element named “amount” defined as PCDATA could be assumed to contain data of type NUMERIC.

Of course, each of these methods can be combined in order to get the best possible assignment of types. As an example of a combined approach consider the following DTD fragment along with a fragment of an XML document conforming to the DTD:

DTD:

```
<!ELEMENT date (#PCDATA) >
```

XML document:

```
<date>01/01/01</date>
```

By applying the sampling method and lexical analysis of element definitions the element `date` could be assigned the type `DATE`. However, even though these methods can be used to provide type bindings for the attributes, complications exist. As described in Section 3, the only data initially fetched from the data source is the DTD. Not knowing the actual structure or size of the data makes it hard to make a sound statement about the statistical value of the sampling. Furthermore, having only a query interface to the data makes it impossible to sample at random, since determining a random position in the XML document is not possible. A way of determining types unambiguously is to read the entire XML document. This is not considered as an option due to the possibly very large documents that must be analyzed. It is beyond the scope of this paper to go further into the problems and solutions associated with type determination. Instead a function `determineType` is assumed, that assigns types from the set *DataTypes* to each of the elements in the set *AttributeList* associated with each of the UML classes.

3 Integration Architecture

This section describes a general system architecture for integrating XML and relational data sources at the conceptual level in a web-based (virtual) DW. The architecture makes the use of XML data transparent to the data analysis tool and thereby to the end user.

XML Data Sources Before the designer can construct a DW based on XML data, the designer needs to know the structure of the data to be used. One way of obtaining this knowledge is to fetch the entire XML document from the Web. This solution is rather unsatisfactory, since fetching the entire XML document can take a considerable amount of time. Another solution is to consider the logical structure of the XML document, described by the DTD. Fetching only the DTD saves time, since the DTD is small in size compared to the document it describes. A potential problem associated with this approach is that the DTD can be overly complex, giving no contribution to the designer's understanding of the structure of the XML data. What is needed is an easy understandable and fast way of communicating the structure of an arbitrary XML data source to the designer.

To accomplish this goal, the DTD is transformed into a Unified Modeling Language (UML) diagram as described in the previous section. Visualizing the structure of XML documents in a graphical way is much easier for the designer to understand than a context free grammar such as the DTD. By describing the structure of XML data sources using a high-level conceptual graphical modeling language, each XML data source on the Web becomes an easy accessible and comprehensible database to both designers and end users. The use of DTDs as a basis for deriving the graphical representation of the data source supports fast browsing of the available data, due to the small size of the DTD. This corresponds to the traditional use of the Web where browsing is probably the most important mean of finding the desired information.

Relational Data Sources An integration of XML and relational data is necessary since an enterprise typically stores different types of transactional data in relational databases. As is the case for XML data, a conceptual model of the structure of the relational data is needed, in order for the designer to make qualified statements about the data to be used. It is not the focus of this paper to describe conceptual models for relational data and it is therefore assumed that UML diagrams describing the relational data can be made available to the designer. Otherwise, database design tools such as ERwin [2] can aid in the process by doing reverse engineering of relational schemas into UML diagrams.

General System Architecture Figure 3 illustrates the architecture which enables the use of XML data and/or relational data in a virtual DW for performing decision support using existing OLAP tools.

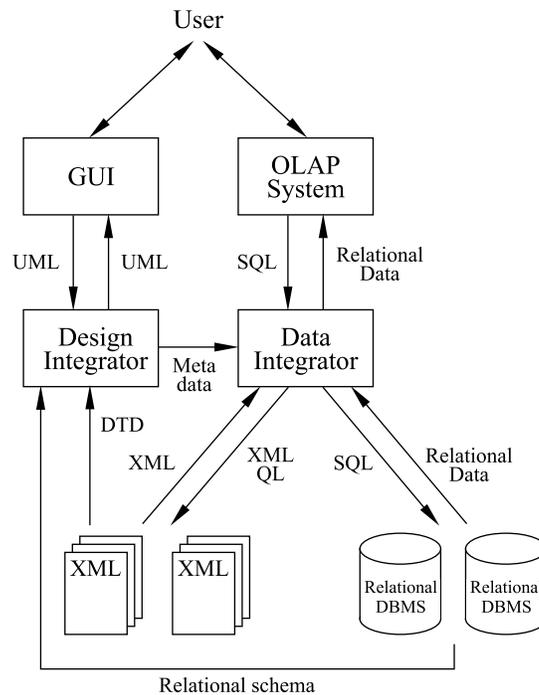


Figure 3: Architecture

The first part of the data integration process is the conceptual integration/design based on the UML diagrams, resulting in a conceptual model of the DW. The DW model (termed a UML snowflake diagram) is constructed by the designer through a graphical user interface (GUI). The UML snowflake diagram is constructed in a semi-automatic way from UML diagrams describing the data sources, which are generated by the Design Integrator. The construction process is described in detail in another paper [8]. The integration/design phase is similar to the “discovery” phase described in the Clío project [12], but is different because we focus on OLAP rather than general data integration and because a conceptual rather than a logical model is used. Finally, the UML snowflake diagram is transformed into relational structures by the Data Integrator, and are then made available to the OLAP System.

At run-time, the data sources are presented to the OLAP tools as relations, creating a relational multidimensional view (a star or snowflake schema [9]) over the data sources. There are several good reasons for using neither UML nor XML as the data model for the run-time system. First, no query language for UML exists, so no query tools would be able to use the constructed DW. Second, almost all data analysis tools, e.g., OLAP tools, can access data using SQL, while only very few analysis tools

can use XML query languages for data access. Third, the distinguishing features of special XML query languages, e.g., querying document order, is not interesting w.r.t. OLAP queries. Fourth, OLAP queries are relatively simple, making the translation from SQL to XML query languages easier. This makes the relational model the obvious choice for the run-time system. The central component in the run-time systems is the Data Integrator which enables both XML data and relational data to be accessed using SQL queries through a standard interface like Open Database Connectivity (ODBC). Integrating XML this way makes the use of XML transparent to the end user since the OLAP tool cannot tell the difference between “true” relational structures and relational structures generated from XML data.

At query time, the user queries the OLAP System using an OLAP query language native to the OLAP System, e.g. MultiDimensional eXpression (MDX) as used by Microsoft’s SQL Server 2000 Analysis Services [11]. These queries are transformed into one or more SQL queries by the OLAP System, used to query the relational structures generated from the UML Snowflake Diagram. The Data Integrator examines the SQL queries delivered by the OLAP System and transforms these queries into a series of queries expressed in either an XML query language or SQL, depending on whether the data originates from an XML source or a relational database management system (RDBMS). The results of these transformed queries (XML from XML data sources and relational tuples from RDBMSs) are then transformed into relational tuples fitting the relational schema generated from the UML snowflake diagram.

4 Conclusions and Future Work

Motivated by the increasing need for data integration and the increased use of XML documents for exchanging information on the Web, this paper considers data integration at the conceptual level, with a focus on OLAP systems.

First, algorithms for automatically constructing UML diagrams from XML data were presented, enabling fast and easy graphical browsing of XML data sources on the web. The algorithms captured important semantic properties of the XML data such as precise cardinalities and aggregation (containment) relationships between the data elements. Second, an architecture for integrating XML data at the conceptual level was presented. The architecture also supported relational data sources, making it well suited for building data warehouses which are based partly on in-house relational data and partly on XML data available on the web. This work is part of a project that investigates technologies for building data warehouses based on web data.

This paper improves on previous work on integration of web-based data by focusing on data integration at the conceptual rather than the logical level. The presented approach also improved over intermediate translations of XML to UML (XML to relational, followed by relational to UML). Another distinguishing feature is the focus on OLAP.

In future work, the first immediate step to be taken is to finish the implementation of a prototype utilizing the aspects described in this paper. A very important aspect of the implementation to investigate efficient query processing techniques such as query translations and data caching. Storing higher-level summaries of the data can speed up query processing considerably. We are currently working on these issues. Furthermore, if XML Schema advances to a W3C Recommendation it would be interesting to consider using this richer formalism for describing XML data sources instead of using DTDs. Other aspects of XML, such as whether preservation of document order is of relevance to OLAP analysis could also be investigated.

References

- [1] Abiteboul, S. et al. Tools for Data Translation and Integration. *IEEE Data Engineering Bulletin* 22(1), pp. 3-8, 1999.
- [2] Computer Associates. ERwin Brochure. www.cai.com/products/alm/erwin/erwin_pd.pdf. Current as of January 18th, 2001.
- [3] Deutsch, A. et al. Storing Semistructured Data with STORED. In Proceedings of *SIGMOD Conference*, pp. 431-442, 1999.
- [4] Fernandez, M. F. et al. Declarative Specification of Web Sites with Strudel. *VLDB Journal* 9(1), pp. 38-55, 2000.
- [5] Florescu D. and D. Kossmann. Storing and Querying XML Data using an RDMBS. *IEEE Data Engineering Bulletin* 22(3), pp. 27-34, 1999.
- [6] Garcia-Molina H. et al. The TSIMMIS Approach to Mediation: Data Models and Languages. *Journal of Intelligent Information Systems* 8(2), pp. 117-132, 1997.
- [7] Hellerstein, J. M. et al. Independent, Open Enterprise Data Integration. *IEEE Data Engineering Bulletin*, 22(1), pp. 43-49, 1999.
- [8] Jensen, M. R., T. H. Møller, and T. B. Pedersen. Specifying OLAP Cubes On XML Data. To appear in *Proceedings of SSDBM 2001*.
- [9] Kimball, R. *The Data Warehouse Toolkit*, Wiley, 1996.
- [10] Lahiri, T et. al. Ozone: Integrating Structured and Semistructured Data, *Proceedings of DBPL*, 1999.
- [11] Microsoft Corporation. Microsoft SQL Server 2000 Analysis Services, www.microsoft.com/sql/product-info/analysisservicesWP.htm. Current as of January 9th, 2001.
- [12] Miller, R. J. et al. The Clio Project: Managing Heterogeneity. *SIGMOD Record* 30(1), 2001.
- [13] Object Management Group, Inc. *OMG Unified Modeling Language Specification 1.3*, www.rational.com/uml/resources/documentation/index.jsp, 1999. Current as of December 13th, 2000.
- [14] Pinnock, J. et. al. *Professional XML*, Wrox Press, 2000.
- [15] Shanmugasundaram, J. et. al. Relational Databases for Querying XML Documents: Limitations and Opportunities. In *Proceedings of VLDB 1999*.
- [16] Roth, M. T. et al. The Garlic Project. In Proceedings of *SIGMOD Conference*, pp. 557, 1996.
- [17] Silicon Integration Initiative *The Electronic Component Information Exchange QuickData Architecture*, www.si2.org/ecix/, 2000. Current as of December 29th, 2000.
- [18] World Wide Web Consortium *Extensible Markup Language (XML) 1.0 (Second Edition)*, W3C Recommendation, www.w3.org/TR/2000/REC-xml-20001006, Oct. 6 2000. Current as of December 20th, 2000.
- [19] World Wide Web Consortium *XML Schema*, W3C Candidate Recommendation, www.w3.org/XML/Schema.html. Current as of December 30th, 2000.