THE FLORIDA STATE UNIVERSITY

COLLEGE OF ART AND SCIENCE

XML AND DATABASE S

By

JUNG KEE KIM

A Survey Paper submitted to the
Department of Computer Science
In partial fulfillment of the
requirements for the degree of
Doctor of Philosophy

Survey Paper Defended in Spring Semester, 2002

The members of the Committee approve the survey paper of Jungkee Kim defended on April, 2002 .

 

_____
Gregory Riccardi
Professor Directing Dissertation

_____
Geoffrey C. Fox
Committee Member

_____
Gordon Erlebacher
Committee Member

_____
David Whalley
Committee Member

_____
Larry Dennis
Outside Committee Member

## TABLEOFCONTENTS

iv

# LIST OF FIGURES

# ABSTRACT

XMLisrapidlybecomingakeytechnologyforinformationrepresentationand
exchangeontheInternet,whereitincreasestheopportunityfortheintegrationofthe
variousdataformats.TheWorldWideWebConsortiumproposesseveralstand          ardsto
definethestructureofXMLandextendtheXMLforaddressingandlinking.Toadopt
suchbenefitsofXML,manyacademicinstitutionsandbusinesscompanieshavemade
XMLenableddatabasesonthecurrentornewarchitectures.Inthispaper,Iwill          present
theXMLstandardsandcomparevarioustechnologiesusedforstoringXMLdataina
repository.Byhighlightingtheirstrengthsandweaknesses,Iseekideasforcombining
thosetechnologieswithacollaboratingsystemtoleveragetheintegrationof          metadata.

# CHAPTER1

# Introduction

Sincet heExtensibleMarkupLanguage(XML)    was introducedas  astandarddocument form,theusageofXMLforthedataexchangeformatbetweenapplications        has dramaticallyincreased .Beforethe   appearanceofXML,applicationd  ependentdata formats wereusedfor   dataexchange .  However,t heWorldWideWeb   environmentisso huge thatapplicationdependentdataformats    requirealotofcoding    effort for synchronizing databetweendifferentapplications  ,duetothe   manydifferentd ataformat ontheWeb .ThoughHyperTextMarkupLangu    age(HTML)protocolisused    onthe Web,itonlydescribeshowthedataareshown        throughWebbrowsers  with the fixedtag format. XMLprovidesuser -definedtagsanditissimple     toproduce  andflexible. With nestedtags,XMLcanrepresentnotonlyhierarchicaltreestructure,butalso        agraph structureusingspecialattributes.  Someapplicationsmayneedanagreement       ontheXML documentformat , becausesuchstandardizationwould   reducetheerrors  andi ncreasethe efficiency inexchangingdata.DocumentTypeDefinition(DTD)andXMLSchema        are usedforforcingthe   productionofXMLdocuments    satisfying particularrules.

Toextractkeyinformationthrought heWeb,metadatacanbeusedtodesignateth       e desiredinformation.TheResourceDescriptionFramework(RDF)provides       a description ofmetadata mainlyrelatingtodataontheWeb      .ThroughRDF,thesemanticmeaningof

theWebresourcescanberevealedtomachine        -orientedsystem s,forexample,search

engines.TheRDFSchemaisdefinedtoconstrain         thevocabularyofRDFdocuments   .

FromtheextensionsofRDFandRDFSchema,asemanticmarkuplanguagecanbe

generatedforthe   intelligentsemanticdescriptionsoftheWebdocuments.Tosupport

linkingf orXMLdocuments,XLink,XPointer,andXPath        wereintroducedbytheWorld

WideWebConsortium.

Anovel      pointofthispaperistosurveytechnologiesforstoringXMLdocumentsina

repository.Manycommercialdatabasesystemcompanieshavetriedtoint          egrateXML

enablingfunctionalitiesontheirexistingarchitecture       sordevisenewarchitecture      sfor

XML document storage.Ifocusonthreetypicalrelationaldatabasecompanies'support

forXMLrepositories  ,and  on otheracademicapproacheswithmapping   X MLdocuments

to arelationaldatabasesystemorasemistructuredrepository.        Thecommercialdat   abase

systemsprovidesomeinternal   orexternal  proceduresandpackagestoshredXML

documentsintorelationalcolumnsorobjects.Thedataintherelationalc        olumnscan

generateXMLdocumentsthroughproceduresandpackagesinthedatabasesystem.

Anotherway  of storing theXMLdocumentsisusingLargeObjecttype      sinrelational

columns.EachLargeObjecttype       recordkeepsthewholeXMLdocument      intheobject

column, withoutshredding.   However,thefunctionalitiesonthosetypes     havenot  yet met

thedemandofXMLusers.     Furthermore,thosemappingsareprogrammedmanually      ,and

someacademicresearch  es approachmapping  ontopofthoseXMLenablingrelational

databasesystems.STORED[14]fromAT&TLabscombinedtherelationalmappingand

semistructuredtechniquestogeneratemapping    sfromXMLdocumentstorelational

tables.ThetechnologyfromUniversityofWisconsin        -Madison[15]usedtheDTDsto

producethe relationalschemaswithseveraldifferent *Inlining*methods.AnotherXML storagetechnologyis thesemistructureddatarepository.    Semistructureddata research startedbeforetheemergenceof   XML,butthesimilarityofsemistructureddataandXML documentsnaturally allowthetechnologytoapplytoXMLdata.LorefromStanford University[2] isasemistructureddatabasesystemsupporting    anXMLrepository. DataGuide[18]isthecorepartforthedescriptionofthestructureofthesemistructured data.

# CHAPTER 2

# XML

TheExtensibleMarkupLanguage(XML)isastandarddocumentspecification,a

subsetofStandardGeneralizedMarkupLanguage(SGML)format[1].ThoughSGMLis

powerfulenoughtohavebeenusedbytheU.S.governmentandpublishingcomp               anies

formakingdocuments,itsimplementationisconsideredtobedifficultandcomplex.The

HypertextMarkupLanguage(HTML)isanotherapplicationofSGML,butHTMLonly

presentstheshapeofthedocumentsontheWeb.ThetagsofHTMLarefixedandth               e

HTMLhasnomechanismforthedatavalidation.ThosefeatureslimitHTMLtoWeb

informationmanipulationandledtotheadventoftheXML.

XMLconsistsofa        *well-formed*structure,rootedinaprolog,oneormoreelements

includingbalancedstart  -an dend -tagswithattributes,andmiscellaneousoptional

featuresincludingcomments,processinginstructions,andwhitespace.Atypical

exampleoftheprologasfollow:

```
<?xml version="1.0" encoding="UTF-8" standalone="yes"?>
```

Thisshowsthatthisisan      XMLdocumentandtheversionnumberis"1.0."Theother

*attributes*,"encoding"and"standalone,"areoptional.Theencodingdenotesthatthis

documentusedUnicodeTransformationFormat8      -bit(UTF -8)encoding.AsXMLis

designedtosupporttheInternation    alcode,theencodingnamescanbeanytheparser

supports, but the name is recommended to be registered with Internet Assigned Numbers Authority. The *standalone* document declaration is also an optional attribute and it indicates whether the XML document is affected by an external markup declaration like DTD and XML schema.

An *element* forms a root of a hierarchical tree structure for an XML document. Other elements can be added and should be nested within each other. If there is no content in an element, the empty-element tag, "<*tag-name*/>," can be used instead of a pair of tags. Those tags are user-defined and that feature makes XML big different from HTML. Attributes are used to attach additional information of an element. They are located in start-tags or empty-element tags. Using attributes of type ID, IDREF or IDREFS[2], XML elements can be uniquely identified and form links. This linking mechanism allows XML to represent graph-structured information as well as tree-structured. In special CDATA sections, any markup data are interpreted as text data.

```xml
<?xml version="1.0" encoding="UTF-8" standalone="yes"?>

<!-- This is an example car. -->
<car id="J544XD" state="NY">
  <company> Toyota </company>
  <model> Corolla </model>
  <type> DX </type>
  <year> 1996 </year>
  <color> white </color>
</car>
```

**Figure 1**. A car description in XML.

Figure 1 shows an example of the XML document. The first line stands for the prolog. The version attribute should be declared and currently the value should be "1.0." The

optional attributes *encoding* and *standalone* show that the code for this XML document is "UTF-8" without any schema ("standalone" is "yes"). The second line is empty -XML allows empty line for the good formatting. The third line is encoded as a comment, which does not present any meaning in the XML structure. The comment only provides optional information to readers. There is only one root element named "car" in the document and other elements are nested in the root element. The car element has two attributes. The *id* attribute identifies a specific car and the *state* attribute means the registered state of the car. The each element opens with start -tag "< *TAG-NAME*>" and closes with end -tag "</ *TAG-NAME*>" and there is no overlapped tags. So, this XML example said to be well -formed.

# CHAPTER 3

## DTDandXMLSchema

DocumentTypeDefinition(DTD)andXMLSchemaarewaystodefinethestructure
ofXMLdocuments.DTDhasbeenusedinSGMLforovertwentyyearsandXMLis
specifiednewlybytheWorldWideWebConsortium(W3C)[3].Th        egoalistomake
rulestoconstructXMLdocuments.Formanypurposes,user        -definedtagsalonedon't
provideasufficientlyrigorousstructurefortheXMLinformationexchanges.By
requiringthesameDTDorXMLSchema,twodifferentapplicationscanagree        ona
particularstructureforanXMLdocument.Ifawell        -formedXMLdocumentsatisfiesa
DTDorXMLSchema,thedocumentsaidtobe        *valid*.

TheXMLSchemaspecificationreflectsthedemandsofusers,whohavefoundDTD
toolimited.Theschemahasman        yimprovedfeaturesoverDTD.BeforedefiningXML
Schema,therewereseveralattemptstoimprovethefunctionalityoftheschemalanguage
forXMLdocuments;someexamplesareDocumentDefinitionMarkupLanguage
(DDML),DocumentContentDescription(DCD),    SchemaforObject -OrientedXML
(SOX),andMicrosoft'sXML   -DataforBizTalk.TheW3Cconsortiumactivityforthe
newschema,XMLSchema,consideredthoseschemasinproducingtheirdesign.The
maindifferencesbetweenDTDandXMLSchemawillbepresented        intheXMLSchema
section.

# DTD

DTDformatisverydifferentfromXML.ADTDisusuallyincludedintheprologpart
ofanXMLdocumentusingthe"!DOCTYPE"tag.TheDTDcanbedefinedexternally
inaseparatefile,designatedwithafilenameoraU        niformResourceIdentifier(URI).
ThetypicalblocksofaDTDare       *elements*and  *attributes*.BNFsyntaxcanbeshownas
follow:

```
<!ELEMENT> <element-name> <element-type>
<!ATTLIST> <attribute-name> <attribute-type> <attribute-option>
```

Figure2sh      owsanexampleDTDforthecardocumentofFigure1.Intheexample,the
"car"elementis  *non-terminal*andtheotherelementsare       *terminal*.Thenon -terminal
element,"car,"hasfivesub    -elements:company,model,type,year,andcolorinthat
order.Iti  scalleda  *sequence*,whichrestrictstheorderofsub        -elementspresent.  *Choice*
isanothergroupoptionforthesub        -elementsanditgivesalistofalternativesforthem.
Theverticalbar("|")isusedasthedelimiterforchoices,andthecommaforseque                nces.

```
<!ELEMENT car(company, model, type?, year, color)>
<!ATTLIST car
          id CDATA #REQUIRED
          state CDATA #IMPLIED>
<!ELEMENT company (#PCDATA)>
<!ELEMENT model (#PCDATA)>
<!ELEMENT type (#PCDATA)>
<!ELEMENT year (#PCDATA)>
<!ELEMENT color (#PCDATA)>
```

**Figure2.** DTDforacardocument

Inthesequenceoftheexample,allbut          *type*elementwillappearexactlyonce.The
typeelementcanbeincludedoptionally.Thisisindicatedbythesuffix,"?."Other

8

allowed suffixes include "+," which means one or more elements can appear, and "*," which means zero or more can appear.

"#PCDATA" in terminal elements stands for parsed character data, which denotes text that has no markup. That is the only way to represent text in DTD, and this was one of motivations for the invention of XML Schema. The element content can also be *empty* or *any*. The *empty* element has no content but may have attributes. The *any* element has no restriction for that element.

The "car" element in the example has two attributes, which are declared in the DTD. The order of the attributes is not constrained. Both of the attributes have the same data type, character data (CDATA). Other attribute types like ID, IDREF, and IDREFS are also very useful and have key roles in representing a graph structures in XML format. The final term is an attribute description specifies whether this attribute is optional or required. The option for the attribute "id" is "#REQUIRED," and this attribute must be appeared in the every defined element. "#IMPLIED" in "state" means the attribute can be optional. Other option is "#FIXED" and this type attribute should have a default value. The fixed value cannot be changeable by the user.

## XMLSchema

While DTD is written in the syntax of Extended Backus Naur Form (EBNF), XML Schema uses XML document syntax. By supporting *namespaces*, XML Schema allows several sources of document definitions to be used in a single document. In DTD, a new DTD is needed to combine multiple DTDs. The XML Schema supports 44 datatypes

including *string*, *decimal*, *time* and *date*, whereas DTD only provides 9 XML    -related

primitive types. Inheritance is another major feature of XML Schema, which is not

present in DTD. This allows reusing existing structures by        extending or restricting the

base types.

Figure 3 includes an example of XML Schema, which reproduces the schema

presented in Figure 2 in DTD format.

```
<xsd:schema xmlns:xsd="http://www.w3.org/2001/XMLSchema">

<xsd:element name="car" type="CarType"/>

<xsd:complexType name="CarType">
  <xsd:sequence>
    <xsd:element name="company" type="xsd:string"/>
    <xsd:element name="model" type="xsd:string"/>
    <xsd:element name="type" type="xsd:string"
      minOccurs="0" maxOccurs="1"/>
    <xsd:element name="year" type="xsd:decimal"/>
    <xsd:element name="color" type="xsd:string"/>
  </xsd:sequence>
  <xsd:attribute name="id" type="xsd:string" use="required"/>
  <xsd:attribute name="state" type="xsd:string"/>
</xsd:complexType>

</xsd:schema>
```

**Figure 3.** XML Schema for a car document

The    *car* schema has one    *schema* element with sub  -elements *element* and  *complexType*.

In the *schema* element, a namespace has been declared. The prefix "xsd:" associated with

the namespace is used on each of the elements. The prefix name of a namespace can be

an arbitrary value and different name    spaces from different sources can be used. In the

multiple namespaces, the different prefix names specify the meanings of elements and

attributes, which are followed by the prefix. In this example, the association forces the

elements and simple types to    be identified with the XML Schema language. In XML

10

Schema, elements may have simple types or complex types. A simple type does not include elements. Many simple types, such as "xsd:string," are defined in the XML Schema. A complex type may have nested elements and carry attributes optionally. The type of the *car* element is defined as a complex type, *CarType*, in the example. As in the DTD example, the elements of a complex type can be ordered with *sequence* tag. All the nested elements except the *year* element have a *string* type. This is declared in the *type* attribute. The *decimal* type in the *year* element is a number. If the year needs to be restricted to four-digit numbers for example, another simple type, *gYear*, already defined in XML Schema, and the declaration can be used instead as follows:

```
<xsd:element name="year" type="xsd:gYear"/>
```

If the *state* attribute has to be two capital letters, a new simple type can be defined as follows:

```
<xsd:simpleType name="StateType">
  <xsd:restriction base="xsd:string">
    <xsd:pattern value="[A-Z]{2}"/>
  </xsd:restriction>
</xsd:simpleType>
```

The equivalent of optional elements in DTD can be expressed in XML Schema using "minOccur" and "maxOccurs" attributes. Additionally to the three restrictions of DTD (*,+,and?), XML Schema can designate any number of minimum and maximum occurrences – for example, between 15 and 30. Any given attribute, the attribute may appear at most once. The "use" attribute designates the attribute usage - one of *required*, *optional*, and *prohibited*.

The *fixed* attribute is used when the allowed value of an element or attribute is unique. For example, I may add another attribute *country* in the "car" element. The state name

11

usedintheattributefromUnitedStates     andthefixedvaluethe    *country*couldbe"US."

Thedefinitionoftheattributeasfollows:

```
<xsd:attribute name="country" type="xsd:string" fixed="US"/>
```

# CHAPTER4

# MetadataSupport

TherehavebeenmanyeffortstoextractusefulinformationthroughtheWeb.
Especially,thesearchenginesdevisedvarioustechnologiestofindtheexactlocationof
thedesiredinformation,butinpracticemanyofthema        lsoshowuselessinformation.
Moreover,theyonlyreachlessthanonepercentofthewholeWeb.Thatisbecausethe
scaleoftheWebissohuge,andkeywordspammingisverywidespread.Thoughsome
Webdirectorysitescate   gorizetheinformationmanually,   t hisisnotforthemachine    -
orientedsystem.

Metadatais"structureddataaboutdata."Thiscouldincludecatalogsoflibraries,
authorlistsofbooks,rankingofWebpagesbyfrequencyofreference,ortherelations
betweenindexes.Bothhumanand       machinegeneratedinformationcanbemetadata.

## RDF

TheResourceDescriptionFramework(RDF)isaW3Crecommendationforastandard
representationofmetadata[5].ThisframeworkisdescribedinXMLformat.RDFhas
aninnatefunctionformachine    -orienteddataexchangingbetweenapplicationsbecauseof
itsXMLfeatures.XMLandRDFprovidesemanticinteroperabilityinthecurrentWeb

domain, but XML only describes the document structure. RDF emphasizes semantic meaning on the Web resources by adding a capability as a data model for knowledge representation.

The basic block of RDF consists of three object types - *resources*, *properties*, and *statements*. A *resource* is anything that can be written as a Uniform Resource Identifier (URI) in the RDF expression. It can be not only a Web page but also an XML element. Anything written in URI could be a resource. A *property* is a specific characteristic, attribute, or relation of the resource - for example, "owner." Each property has a specific meaning, which can be classified by a schema related to the name of the property. A *statement* is a combination of a *resource*, a *property*, and a *value*. Each part of a statement is also known as the *subject*, the *predicate*, and the *object*. The object can be another resource or a literal, which might be any string or XML. Figure 4 presents an example of RDF graph for the Web page of the University.
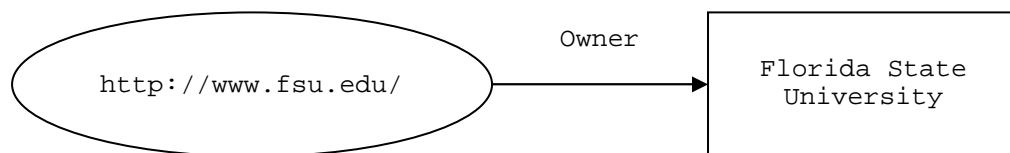


**Figure 4**. An example of RDF graph

In the figure, the oval shape node denotes a subject, the arc denotes a named property, and the rectangular shape is a node, which represents a literal. The graph represents the following statement:

```
    "Florida State University is the owner of the resource
http://www.fsu.edu/."
```

It also can be read as:

14

```
          "http://www.fsu.edu/ has owner Florida State University."
```

ThestatementcanbewrittenintheXMLformat:

```
<?xml version="1.0">
<rdf:RDF
  xmlns:rdf="http://www.w3.org/1999/02/22-rdf-synfax-ns#"
  xmlns:s="http://description.org/schema/">
  <rdf:Description about="http://www.fsu.edu/">
    <s:owner>Florida State University</s:owner>
  </rdf:Description>
</rdf:RDF>
```

TheRDFXMLsyntaxhasarootelement,<  *RDF*>,butthiselementisoptionalwhenthe

*description*isknowntobeRDFfromtheappl    icationcontent.InRDFelement,the

namespaceattributesdesignatethelocationofthedeclarationsoftheRDFelementswith

theprefix"rdf:",andthelocationoftheschemadeclarationassociatingwiththeprefix

"s:".Thenamespacedeclarationcanal    ternativelyappearinaspecific  *description*

element,oreveninpropertyelements.The       *description*elementhasthesubject;thechild

elementsdescribethepropertiesandtheobjects.Inthisexample,"<s:owner>"and

"</s:owner>" -apairoftags   -show theproperty.Theobjectis"FloridaState

University."

AsinXMLSchemaforXMLdocument,RDFSchemaprovidesavocabulary

constraintfacilityforRDFdocument.InRDFSchema,theclassesoftheresourcesare

defined.Theclasseshavethesamerole    asintheobject   -orientedprogrammingmodels.

Theclasseshavehierarchicalstructuresandtheyareextendedwithsubclassrefinement.

Thetermssuchas"  *Class*,""  *subPropertyOf*,"and"  *subClassOf*"areusedforthebasic

typesystemforRDFtodefinesuch     classes.Byusingclassconcepts,thereusabilityof

metadatacanbeincreasedbecausesharingschemasandaddingsubclassestotheexisting

schemaswillproducesufficientmechanismsinmanyschemaspecifications.

Ontologies play a crucial role in t he "Semantic Web" – a machine - understandable Web with intelligent services. They provide  shared and  precisely defined  terms in a particular domain for communications between human users and application systems. DAML+OIL – combined terminology from old ve rsions, DARPA Agent Markup Language (DAML) and Ontology Inference Layer  (OIL) - is an ontology language submitted to  W3C as a semantic markup language for Web resources [19].  It extends RDF and RDF schema. The object -oriented structure of domain s in DAML +OIL consists of the terms " *Class* " and " *Property.* " The  usage of  the term  class is similar to that of RDF schema, but  the classes of DAML+OIL are less restricted. For example, "  *subClassOf* " class elements of DAML+OIL allow cyclic subclass -relations. A pro pertyis a binary relation, which defines the  relation between two items. All expressions of description logic can be written in DAML+OIL terms.  The DAML+OIL properties  are of two types: "*objectProperty*" for  object relations and " *datatypeProperty*" for da tatype values.  Those properties are ranged and m ultiple ranges  can apply to a property  *conjunctively*. For example, positive integer range and greater than 1 range will produce a property that has positive integers greater than 1.  The restriction of  domains in DAML+OIL is  global, while the  RDF Schema  only has a single range and a local scope on domains. However, efficient tools for  reasoning about Web resources and  complete algorithms for the full DAML+OIL language are not provided yet  [20] .

From th eviewpoint of semantic interoperability, RDF is better than  plain XML. RDF vocabularies are simple enough to manipulate huge number sof data. Meanwhile, XML often regards ordered elements as important and has complex structure. Those features make it d ifficult for XML to handle large amount of data. Additional data conversion is

notnecessaryinRDFbecauseRDFpresentsdomainmodelsnaturallywithdefining objectsandrelations.AnotherbenefitofRDFisindependency oftheXML.In an XML document, aschemachangemaycauseinvalidityforthequerybasedon theoldstructure ofthe XMLdocument .RDFpresents a semantictree that isparsedwithonlyusableset oftriples ,andthedatanottobeinterpretedareignored[6].

## XLink,XPointer,andXPath

AsinhyperlinksofHTMLdocuments,XMLdocumentscanbelinkedtootherXML documentsbyusingXMLLinkingLanguage(XLink)[7],XMLPointerLanguage (XPointer)[8],andXMLPathLanguage(XPath)[9].TheXLinkdescribeslinks betweenresources.T heXPointerpointsthereferencethroughURI.TheXPathpresents thelocationofspecificpartsofanXMLdocument. Lastversion oftheXPointer (September2001)isaCandidateRecommendation oftheWorldWideWebConsortium. The XLink andtheXPathisr ecommendedbythe W3C.

**XPath**

ThemainpurposeofXPathistodesignatepartsofanXMLdocument. XPath providesanextendedaddressingsyntaxthatdefinesacompactnotationfor *node*location intheXMLdocumenttree.XPathdoesnotuseXMLsyntaxbu titisastring -based language.TheExtensibleStylesheetLanguageTransformations(XSLT)andthe XPointerusethefunctionalityoftheXPath.

**XPointer**

XPointerisusedtoidentifyspecificfragmentsinXMLdocumentsviaaURI.

XPointermayselecto nthebasisofXMLIDattributes,ornodesinthehierarchical

structureofanXMLdocumentrelatedusingXPath.AnXPointercanalsoreferencean

arbitraryuserdesignationonaspecific *point*or *range* –doesn'thavetobeanXMLnode.

The *range*canbe specifiedwithtwo *points*,as:

```
xpointer(id("start")/range-to(id("end")))
```

ThatXPointerlocatestherangebetweenthestartpointfortheelementwithID" *start*"

andtheendpointfortheelementwithID" *end*".TheXPointeriscomplexenoughto

presentmostusages.

**XLink**

XLinkisabletolinknotonlydocumentsbutalsoresources,whichincludedocuments,

audio,video,databasedata,andanyaddressableinformationorservices.WhileHTML

linksneedtoedittheresourceforadditionallinks,XLink don'trequireanywrite

permissiontoeditthesource.TheXLinkcansimplysettheURIwiththestartingand

endingpointforthelinking.TheXLinkalsoprovidesmultidirectionallinks( *extended*

links)aswellastheunidirectionallink( *simple*link) –thetraditionallinkontheWeb.

Thelinkscanbestoredexternally( *extended*link)ofthedocuments,theyaddresswith

URI,andtheycanbeinline.Traversalof" *A*" linkusuallyreplacesthedocument

currentlyviewed.Traversalmeans"usingorfollowin galinkforanypurpose"[7].The

usermayinitiatetraversalwithclickingonthelinks,ortheretrievingdocumentmay

initiateit.

```
<?xml version="1.0"?>

<doc xmlns:xlink="http://www.w3.org/1999/xlink">

<head>
<title> Animals </title>
<extendedlink xlink:type="extended">
  <loc xlink:type="locator"
       xlink:label="seaplace"
       xlink:href="#xpointer(//body/animal[1]/sentence[2]/place[1])"/>
  <loc xlink:type="locator"
       xlink:label="seareference"
       xlink:href="#sea"/>
  <arc xlink:type="arc"
       xlink:from="seaplace"
       xlink:to="seareference"
       xlink:show="new"
       xlink:actuate="onRequest"/>
</extendedlink>
</head>

<body>

<animal name="whale">
<sentence>Whales are mammals.</sentence>
<sentence>Whales live in the
  <place>sea</place>.</sentence>
</animal>

<animal name="horse">
<sentence>Horses are mammals.</sentence>
<sentence>Horses live in the
  <place>land</place>.</sentence>
</animal>

</body>

<tail>
<reference>
<places id="sea">
  Sea is the continuous body of salt water covering the earth.
</places>
<places id="land">
  Land is the part of the earth not covered by water.
</places>
</reference>
</tail>

</doc>
```

**Figure5** .AnexamplefortheXMLdocumentwithXLinkexpressions.

19

Figure5showsanexample,whichincludesthefeaturesofXLinkandXPointer.Inthe

example,"`<doc xmlns:xlink = "http://www.w3.org/1999/xlink">`"denotesthe

XLinknamespacedefinitionwiththeURI.        The *extendedlink*elementisakindof

extendedlink,whichhasfullXLinkfunctionalitysuchasarcs(        *inbound*and  *third-party*)

andlinkswitharbitraryresources.Theoth        ertypeforthelinkisthe        *simple*link,which

hasonlytwoparticipatingresources.

Inthe        *extendedlink*element,threesub  -elementsareem  bedded:theyaretwo  *locator*

elements(XLinktype)andan        *arc*element(XLinktype).The        *locator*typeelement

designatesremoteresources,whoselocationisdenotedwiththelocatorattribute,"        *href*".

Thearctypeelementsrepresentthelinktraversal,wh        ichisusuallyapairofstart(  *from*)

andend( *to*)resources.The  *locator*labeled" *seaplace*"hastheXPointerandthe

expressionintheparenthesesofthisXPointerisalsotheXPathexpression.This

XPointerpointstothefirst        *place*elementoftheseco  nd *sentence*inthefirst  *animal*

elementofthe  *body*element.The"  *seareference*" *locator*linkstothefirst"  *places*"

element,whichhastheattribute        *id*named" *sea*."

The        *arc*elementintheexamplehastworemoteresourcesforthetraverseanditis

calleda  *third-party*arc.Ifthearcfromlocalresourcetoremoteresource,itisthe

*outbound*arc.Or,the        *inbound*arctraversesfromremoteresourcetolocalresource.The

traversalattributes," *from*"and" *to,*"areforthestartandtheendpointsofth        elink.The

*show*and  *actuate*attributesrepresentthebehavioroftthelink.Theydesignatethe

behavioroftheendingresourceofthearc.Inthefigure,the"new"valueforthe        *show*

attributewillopenanewwindowwhenthetraversaleventhasbeenreq        uested.The

*actuate*attributesets" *onRequest*,"anditconstraintsthetraversalevent.Ifthevalueof

20

the *actuate* attribute changed to " *onLoad*," the new window would be shown immediately

on loading the starting resource.

# CHAPTER5

# XMLandDatabases

AnXMLdocumentisatextdocumentusedforinformation exchange between applicationprograms , typically throughtheWeb.XMLdoesnotforceanyinternal structure onthecomputer. However,itisimportanthowtostore andprocesstheXML document fromtheviewpointofefficiencyofdatamanipulati on.XMLlogicallyhas a treestructurewithelementsandattributes ,and isused bothin theroleof a data transport formatand as a document markuplanguage.Bourret [10]c lassifiestheXMLdocuments as *data-centric*documentsand *document-centric*documents. Data-centricdocuments are highlystructured and haverelativelysmallsizedtext elements.Document -centric documents havefree -formattext withsomewordsmarked -up. Theperform anceofthe XMLdocument processingdependsonthe kind ofXMLdocumentpresentation:data - centricanddocument -centricdocumentshaveprosandcons relativeto thedifferenttypes ofdatabase .

Manyattemptsweremadetoleveragethecurren tdatabasetechnologies torepresent XMLdocuments .Mapping XML torelationaldatabasewasoneoffirstmethods used to storeXMLdocument s in existingdatabase s.ManyXMLrepositories haveproposed waysto map XML torelationaldatabase s, becausethe re lationaldatabasedominatesthe currentdatabasemarketandmanyapplications arealreadydevelopedonrelational

databases. Notonlylegacyapplications,butalsodata -centricXMLdocumentsobtain benefits fromusing therelationaldatabase. Object-relational storageis anatural fitfor XMLstorage ,because itslogicalstructure is similarto thatofXMLdocuments. Another approachfortheXMLstorage is from *semi-structured*data. Semi-structureddatahad beendevelopedbeforeXMLstandardemerged.Th eXML format unsurprisingly applied to thesemi -structureddatabase ,because XMLisasemi -structureddataformat. Lore[ 2] isawell -knownexampleofsemi -structureddatasystemtopresentXMLdocumentsas semi-structuredata.

## RelationalDatabase

### OracleDatabaseSystem

Currentlytherelationaldatabasedominatesthedatabasemarket.SinceXML has emergedasanewstandardfortheinformationexchange,manyrelationaldatabaseshave beentriedtocombinetheirdatabasesandXMLtechnologies.Oneof initialanswers fromthecommercialdatabasesistheXML -enabledOracle8ifromOracle, whichhasthe biggestmarketshareinthecurrentdatabasebusiness.

MappinganXMLdocumenttoatableorseveraltablesisa primitiveway ofstoring XMLdocumen tsontherelationaldatabase. Theelements,attributes,andnamesare mappedtothecolumnsoftables inwaysdependingonthefunctionalitiesofdatausage andmappingdesign. Thismethodisusefulintransferringdatabetweenrelationa l databases,but itisnotapplicabletosophisticatedXMLformats.

*Oracle*[11] providesanaturalmannerusingobject    -relationalsupport.   The elements
withattributesdefine  *object*types  thatencapsulatedata.    Setsofobjecttypesand
referencestoobjecttypes   ca nform amodelof    classes.Aclassmapsintoatable.     In
Figure6,  *CAR_TYPE*isdefinedasanobjectandasingleobjectmapstoatableinthis
example.

```
create table CARS
(
ID          VARCHAR2(7),
STATE       VARCHAR2(2),
COMPANY     VARCHAR2(16),
MODEL       VARCHAR2(16),
TYPE        VARCHAR2(3),
YEAR        NUMBER(4),
COLOR       VARCHAR2(16)
)
/

create or replace view NEWCARS as
select SYS.XMLTYPE.CREATEXML('<CAR/>') "CAR"
from CARS;

create or replace trigger CAREXPLOSION
instead of insert on NEWCARS
for each row
declare
  CARID       VARCHAR2(7);
  STATE       VARCHAR(2);
  COMPANY     VARCHAR2(16);
  MODEL       VARCHAR2(16);
  CARTYPE     VARCHAR2(3);
  CARYEAR     NUMBER(4);
  COLOR       VARCHAR2(16);

  DOCUMENT    sys.XMLTYPE;
  ELEMENT     sys.XMLTYPE;

  NOT_A_CAR   exception;

  I           binary_integer;

begin

  DOCUMENT := :new.CAR;

  if (DOCUMENT.existsNode('/car') = 0) then
    raise NOT_A_CAR;
  end if;

  CARID := DOCUMENT.extract('/car/@id').getStringVal();
```

24

```
   STATE := DOCUMENT.extract('/car/@state').getStringVal();

  ELEMENT := DOCUMENT.extract('/car/company/text()');
  if (ELEMENT is not null) then
    COMPANY := ELEMENT.getStringVal();
  end if;

  ELEMENT := DOCUMENT.extract('/car/model/text()');
  if (ELEMENT is not null) then
    MODEL := ELEMENT.getStringVal();
  end if;

  ELEMENT := DOCUMENT.extract('/car/type/text()');
  if (ELEMENT is not null) then
    CARTYPE := ELEMENT.getStringVal();
  end if;

  ELEMENT := DOCUMENT.extract('/car/year/text()');
  if (ELEMENT is not null) then
    CARYEAR := ELEMENT.getNumberVal();
  end if;

  ELEMENT := DOCUMENT.extract('/car/color/text()');
  if (ELEMENT is not null) then
    COLOR := ELEMENT.getStringVal();
  end if;

  insert into CARS values (CARID, STATE, COMPANY, MODEL, CARTYPE,
CARYEAR, COLOR);

  exception

    when NOT_A_CAR then
      raise_application_error(-20000, 'Only car documents can be
stored in this column.');

end CAREXPLOSION;
/


create or replace type CAR_TYPE as object
(
"@id"            VARCHAR2(7),
"@state"         VARCHAR2(2),
"company"        VARCHAR2(16),
"model"          VARCHAR2(16),
"type"           VARCHAR2(3),
"year"           NUMBER(4),
"color"          VARCHAR2(16)
)
/

create or replace view CARDOCUMENTS as
select
sys_xmlgen(
CAR_TYPE(
```

```
  C.ID,
  C.STATE,
  C.COMPANY,
  C.MODEL,
  C.TYPE,
  C.YEAR,
  C.COLOR
),
sys.xmlgenformattype.createformat('car')
) car
from CARS C;


SQL> select C.car.getClobVal() "car" from CARDOCUMENTS C;

car
---------------------------------------------------------
<?xml version="1.0"?>
<car id="J544XD" state="NY">
 <company>Toyota </company>
 <model>Corolla </model>
 <type>DX </type>
 <year>1996</year>
 <color>white </color>
</car>
```

**Figure6** .AnexamplefortheXMLdocumentw      ith mappingtoOracledatabase   .

AnotherapproachforXMLdocumentstorageis      *LargeObject* (LOB).ALOBcolumn

holdsanXMLinstanceand      thistypeofstorageisusefulfor      document-centric

documents.TheXMLinstanceinLOBcolumnscanbeindexedasino            thertexts. The

Newversion of  the Oracledatabase , *Oracle9i*  database system provides *XMLType* - an

objectdatatypeon    a *characterlargeobject* (CLOB)column  storage .XMLType

supportsXPath  inSQLqueries  toextractelementsandattributes      ofXMLinstan  cesas

follow:

```
  SELECT c.car.extract('/car/model/text()').getStringVal() FROM cars c;
```

   ToimproveperformanceakeyelementoftheXMLdocumentcanbestored            in

anothercolumnandbe    indexed.Withoutseparatecolumnindexing,      extracting query

performanceis notacceptable inthecaseoflarge numberof rows. Updatingelementsor attributesinLOBisnotpossible. OnlytheentireXMLdocumentsupdateisallowed. Therearea dditionalpackagesandfunctionsintegratedinSQL whichcanbeusedin queriestowrapthedataincolumns andproduceXMLdocuments. Figure6showsanexample thatmap sbetweenanXMLdocumentandarelational tableusingOracle9i .Thefirsthalf ofthecode representsthemappingfromanXML documenttoarelationaltable .ThelasthalfshowshowtogetanXMLinstancefroma relationaltable.The *CARS*tablestorestheelementsandtheattributesof theXML examplefromchapterone . The *NEWCARS*viewshowsXMLinstancesfromarelational table, *CARS*.WhenanXMLdocumen tinsertsintoarelationaltable,eachelementand attributeneedstobeshred ded.The *CAREXPLOSION*triggersubstitutestheinsertquery tomaptheXMLdocumentintopropertypesofcolumns.Fromtherelationaltable, an XMLinstancecanbeproduced.T he *CARDOCUMENT*viewgeneratesanXMLinstance fromatableusingSYS_XMLGENfunctionintheSQLquery.

**IBMDB2DatabaseSystem**

IBM's *DB2*databasesystem [12] alsosupportsXMLdocumentsforstorageandquery. TostoreXMLdocumentsin *IBMDB2XMLExte nder*, anXMLrepository, *XMLcolumn* and *XMLcollection* optionsareavailable .

In the XMLcolumnoption,theXMLdocumentissavedas *XMLCLOB*, *XMLVARCHAR*or *XMLFILE*typewithout shredding.Someelementsandattributescan bestoredin *sidetables* and i ndexedforperformanceimprovement .Thesidetablesand indexmustbe in a *DataAccessDefinition* (DAD). An exampleDADforsidetablesis given inF igure7.

In the XMLcollectionoption, anXMLdocumentmapsintoasetofrelationaltables andthema ppingmechanismsbetweenDTDandtablesaredescribedin a DAD.This optionisusefulforfrequentupdate s topartofdata,extraction ofonlysomepartofdata, and relatingtootherrelationaldata.

```
…
<dad>
  <dtdid>dxx_install/dtd/getstart.dtd</dtdid>
  <validation>YES</validation>
  <Xcolumn>
    <table name="car_table"> </table>
      <column name="year" type="decimal(4,0)"
       path="/car/year"/>
    </table>
  </Xcolumn>
</dad>
```

**Figure7** .AnexampleofDADsidetabledefinit ion

TopublishXMLdocument sfromdatabasetables, SQL queries with macrosofascript languageandstoredproceduresareused.

**MicrosoftSQLServer**

*MicrosoftSQLServer* [13] isanotherrelationaldatabasesystemthatgeneratesand storesXMLdoc umentsthroughrelationaldata.XMLenablingfeat uresareexecutedin middletier applications ,forexample, *templates*and *XMLviews* . TemplatesareXML documentswhichincludeSQLqueries executedagainstthedatabase.XPathissupported intemplates. An XMLview oftherelationaldatacanbecreatedbyannotatedschema using *XMLDataReduced* (XDR)schema. Theseannotationsareusedtospecify a mappingbetweenXMLandrelationaltables.

Topublish an XMLdocumentfrom a relationaldatabase,SQLSe rver includesSQL extensionstoproducequeryresult sasXMLdocuments.Therearethreedifferentways

toserializeSQLqueryresultsin toXML: *RAW*, *AUTO*, and *EXPLICIT* mode s. InRAW

mode,eachrowofthequeryresultmapsintoa nameofan XMLelement, *row*,andeach

non-NULLcolumnofthequeryresultmapstoanXMLattribute. The examplequery

mayproduceasfollow:

```
SELECT CarID, CarState FROM Cars FOR XML raw

<row CarID="J544XD" CarState="NY"/>
```

Toproducequeryresult s withnestedXMLelements,SQLse rverprovidesAUTOmode.

EachrowmapstoanXMLelementandeachtable alias isusedfor theelementname.

Theorderofthetablenames intheSELECTclause determinesthenesting accordingto

their appearancefrom lefttoright.Thecolumnsoftheque ryresultsmaptotheXML

attributesasinRAWmode. Theexampleofanautomodequerymayreturntheresultas

follow:

```
SELECT Dealers.DealerID, Dealers.DealerName, Customers.CustomerName
FROM Dealers, Customers
WHERE Dealers.DealerID = Customer.DealerID
FOR XML auto

<Dealers DealerID="ROMANOTOYOTA" DealerName="Peter Dan">
  <Customers CustomerName="Jungkee Kim"/>
  <Customers CustomerName="Bryan Carpenter"/>
  …
</Dealers>
```

TheEXPLICITmodecanproduceanyXMLdocumentfromtherelationaldatabase

tables. TheEXPLICITmodequery , whichspecifies thestructureoftheXMLtree,

producesa *universaltable* whichconsistsof *tag*and *parent*,columnnames,androw

ordering. For,examplethequerymayproduce theuniversaltableandanXMLinstance

asfollows:

```
SELECT 1 as Tag, NULL as Parent, Dealers.DealerID AS
[Dealers!1!DelaerID],
  NULL as [Customers!2!name!element]
FROM Dealers
UNION ALL
```

```
SELECT 2, 1, Dealers.DealerID, Customers.CustomerName
FROM Dealers INNER JOIN Customers ON Dealers.DealerID =
Customers.DealerID
ORDER BY [Dealers!1!DealerID]
For XML explicit

Tag Parent Dealers!1!DealerID Customers!2!name!element
1   0      ROMANOTOYOTA        NULL
2   1      ROMANOTOYOTA        Jungkee Kim
1   0      SYRACUSEFORD        NULL
2   1      SYRACUSEFORD        Bryan Carpenter

<Dealers DealerID="ROMANOTOYOTA">
  <Customers><name>Jungkee Kim</name></Customers>
</Dealers>
<Dealers DealerID="SYRACUSEFORD">
  <Customers><name>Bryan Carpenter</name></Customers>
</Dealers>
```

Toproduceaconvenientrelationalviewfrom    an XMLdocument,  *OpenXMLrowset*

*provider*isprovided ,similartothe   *extract*functionofOracle9i.

**STORED**

STORED(SemistructuredTORelationalData)[14]isoneof        several initialproposals

forstoringandqueryingXMLdocuments     mappedtorelational  databasesystem.

STORED utilizes a combinationofrelationalandsemistructuredtechniquestomanage

XMLdocuments.  It triestodiscover   the mostfrequent lyoccurring sub -treesfrom  XML

documentsandmapstheextractedsub    -treestorelationaltables.   Theremaining partof

XMLdocuments isstoredinasemistructured    *overflow*graph.  Forexample,theXML

document andthegraphinF   igure 8 willproducearelational   table *car*with  attributes,

plateanddriver. The   seconddriver,Mary,  mightbestoredin    the overflowrepository .

Adeclarative   query languageisdefined   andthislanguageexpressesthestructurein

bothinputandoutputofaquery.    Becausethe  languageisnon  -recursive,theelements

cannothavearbitrarynumberofsub    -elementsinthesameformat.If     a DTDisavailable,

theoverflowqueriesareexpectedtobe     executedfasterthan  theywouldbe  withouta

DTD.

```
<cars>
  <car>
    <plate>J544XD</plate>
    <driver>Jake</driver>
  </car>
  <car>
    <plate>S7RDG</plate>
    <driver>Tom</driver>
    <driver>Mary</driver>
  </car>
  <car>
    <plate>53G123</plate>
    <driver>Bryan</driver>
  </car>
</cars>
```
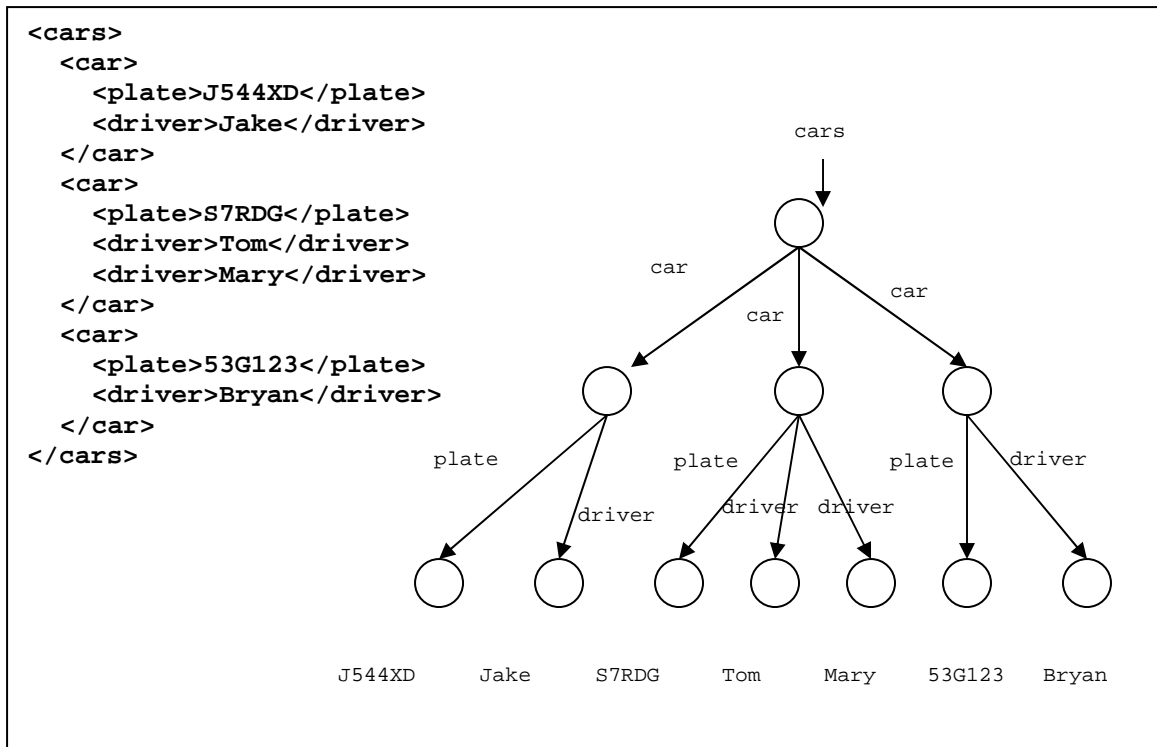
**Figure8** .AnExampleofanXMLdocumentandagraph

WhenSTOREDgenerates     a mapping, itrequires  various parameters : the maximum

numberoftables,  the maximumnumberofattributespertable,      the maximumdiskspace,

the *collectionsizethreshold* ,and  *minimumsupport* .  Thec ollectionsizethreshold

designatestheboundarybetween    *smallsets* and  *collections*.  Them inimumsupport

parameter isrequire dfortheminingalgorithm.Theminingalgorithmdetermineswhich

elementsoftheXMLdocumentshouldbestoredinrelationaltable       sandwhichelements

shouldbesavedinoverflow.     Them iningpro cessesmeasurewhichpathprefixesand

bodiesare  frequentlyoccurring inXMLdocumentsandqueries.Fromthosenumbers,

thebest relationalmappingwillbeselected     takingintoaccount  thenumberofmatching

valuesand  the maximumdiskspace.

31

Sincet hequeriesandupdate areperformed againsttheoriginalXMLdocumentsrather thanmappedtablesoroverflow,STORE D should *rewrite* themintoqueriesandupdates over relationaltablesandoverflowstorage. The rewritealgorithm of the STORED system, *inversionrules* , convertsthequeries overXMLdocuments totake table columns andadd treestructuretothedatalayout.

**RelationalMappingTechnologyfromUniversityofWisconsin -Madison**

Shanmugasundaramandothers [15] suggestedadifferentapproachfor mapping XML documentstorelationaldatabasesystem s.Theyuse a DTDtogeneratearelational schemaunlike STORED,where using a DTDisoptional. Thereisnoconsideration of thequeryworkload. XMLdocumentsareparsed, matchedtoDTDs,andloadedto relationaldatabasetables.Theyused an IBMDB2database fortherelationaltables . For queryingthedata,semistructuredqueriesaretranslatedtoSQLqueriesandtheresultsare convertedtoXML.

WhengeneratingrelationalschemasfromDTDs,DTDs aresimplifiedwithasetof transformations thatismorerestrictedthan thatofSTORED. Transformationsare performedforflattening ,simplification,andgroupingas inthe follow ingexamples:

```
Flattening: (e₁, e₂)* -> e₁*, e₂*
Simplification: e₁** -> e₁*
Grouping: …, a*, …, a*, … -> a*, …
```

Becauseofthesimplification,theorderofelemen tsoftheoriginalXMLdocument may belostthough additionalfieldsforsomeelements weresuggestedforkeepingtheorder. ADTDcanbeexpressedasagraph. Thes implifiedDTDgraphscanbeconvertedto relationalschemasusingoneofthreeproposedmethods:the *BasicInliningTechnique* , the *SharedInliningTechnique* ,orthe *HybridsInliningTechnique* .

The BasicInliningTechniquecreatesseparatedrelationsfo reachelementsofanXML

document,becauseallelementsofaDTDcanbearootofanXMLdocument.Each

relationhasanID thatactsas the keyfield.Allsub -elementsandattributesofthe

elementintherelationareinlined,buttherearetwoexceptio ns.M ultiplyoccurringsub -

elements –thenodesbelow"*" –andrecursivereferencedelementsarenotinlined.

Theyhaveseparaterelations. Therecursiveelementwillalsohave theparentreference.

TheXMLfragment inF igure 8mayproducethe relatio ns as following:

```
cars(carsID: integer)
cars.car(cars.car.carID: integer, cars.car.parentID: integer,
cars.car.plate: string)
cars.car.plate(cars.car.plate.plateID: integer,
cars.car.plate.parentID: integer, cars.car.plate: string)
cars.car.driver(cars.car.driver.driverID: integer,
cars.car.driver.parentID:integer, cars.car.driver: string)
```

The SharedInliningTechnique expresseseachelementasarelationandavoidsthe

redundancyoftheBasicInliningTechnique.Themultiplerelations forelementswit h

severalparents inBasicInliningarestoredineachrelation inSharedInlining. Anew

relationforthoseelements iscreatedand shared.Allthenodeswithin -degreeofonewill

beinli nedintocolumnsofthetablesfortheirparentelements. TheS hared Inliningfor

Figure8mayproduce the followingrelations:

```
cars(carsID: integer, cars.isRoot: boolean)
car(cars.car.carID: integer, cars.car.parentID: integer,
cars.car.parentCODE: integer, cars.car.plate: string, cars.car.plate:
string)
driver(driverID: integer, parentID:integer, cars.car.driver: string)
```

However,theSharedIn lining canrequiremorejoinoperationsonsomeparticular

elementscomparingtotheBasicInlining.

TheHybridInliningTechniqueissimilartotheSharedInliningTechnique ,butithas

additionalinliningwhichisnotincludedintheSharedInlining.Theextrainlining

elementsinthistechniquearenot forrecursiveorcyclic.

33

TheperformancefortheBasicInliningTechniqueisverypoor.Therearetrade            -offsfor
theS haredInliningTechniqueandtheHybridInliningTechnique.


## SemistructuredDatabase


Research on semistructuredd ata initiallylooksfor   anefficientwaytodescrib   e data
withno  fixedschema,forexample,the     informationontheWorld   -WideWeb.  The
relationalorobject  -oriented databasesystem sha veschema ,and  allthe  instancesinthe
systemsshould  becreated  under therule sof  the schema.  Semistructured dataha veno
explicitexplanationofthestructure    , butthey  includedirect  descriptions of thedata  witha
simplesyntax.  TheObjectExchangeModel(OEM)    isatyp  icalsemistructureddata
model.  It wasdesignedforexch  angingdatabetweenheterogeneoussystems    and
originatesfrom  the Tsimmis[ 16]dataint  egrationproject.   AnOEM  objecthasfour
components: *label*, *oid*, *type*,and  *value*.  Labelisacharacterstring    and  oidisthe
identifier oftheobject.   Typecanbeanato   mictype  or a complextype.Ifthe   object type
iscomplex,  the value is asetofoid   s.  Theva lueofanatomictypeis     oneofbasetypes   –
*integer*, *string*, *image*, *sound*,etc .  ThereforeOEMdataformagraphinwhichthenodes
aretheobjectsandthe    labelsare  generally attachedto  edges,thoughtheinitialsystem
was labeled onnodes .

Abiteboulandothers  [ 17] announcedth ree approachestodevelopadatabase
managementsystemforsemistructureddata:    *buildinganapplicationontopofexisting*
*relationalorobjectdatabasesystems*    , *using alow -levelobjectserver* ,and  *buildingthe*
*systemfromscratch* . T heLore(Light  weightObjectReposito  ry)system fromStanford

University[2] tookthelastapproachand is atypicaldatabasesystemforsemistructured data. The Lore system consistsof the API, queryprocessors,and data managers. A graphicaluserinterfacesprovid es queriesandviewsofdata bytheusers . *Lorel*is the querylanguageofLore . Thequeries presented withLorelare parsed, preprocessed, transformed, andoptimize dinqueryprocessors . Thereare several managersandtools fordata handling.TheOEMobject s canbesavedto orfetchedfromfilesthough an object manager. TheXMLdocumentsare mappedand keptinOEMform s internally. Todescribe thestructuresofthe storedsemi structureddata, DataGuides[ 18]were introducedinLore .

ADataGuideisa *concise*and *accurate*summaryof thestructureof a specifieddata graph. I t describeseveryuniquepathof theoriginaldatagraphonlyonce . Everylabel path appearingin theDataGuide existsintheoriginal graph to reservetheaccuracy. The logical structureofaDataGuideis adirected acyclic graph. A *targetset* is thesetofall objects thatagive nlabelpathinagraph reaches. Anodeinthe inputgraphmayappear morethanonce inthetargetset ,becausethenodecanbe reached from several different edges. Thetargetsets containing IDs are the *annotations*of the labelpathsthatdesignate thenodes oftheDataGuide. Goldman andWidom[18] assert that generating DataGuides overtheinputgraphs issimilarto the conversionfroma non-deterministicfinite automation(NFA)to adeterministicfiniteau tomation(DFA). Theyexpect that the conversion wouldbefinishedwithina reasonable timeandtheydi dnotsee any exponentialti meorspace problem duringtheirexperiments. However,thoseannotations mean that multiplelabelpaths mayexisttoreachthesameobject. In Figure9 ,a data graphandtw oDataGuidesofthegraph a reshown .

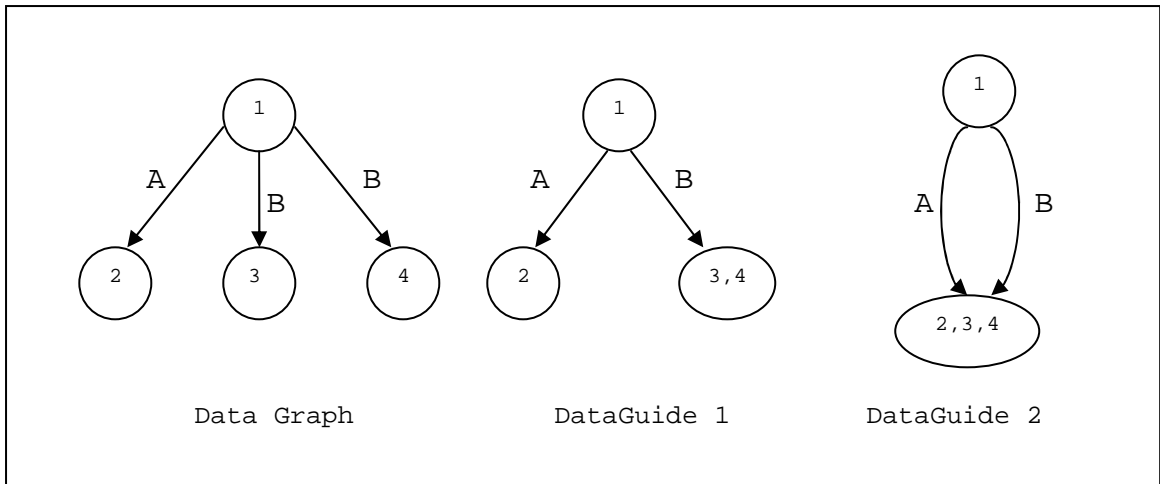**Figure**9.AdatagraphandtwoDataGuides

Wecan reach thenode2 withpath *A* from thedatagraphandDataGuide1 ,but thesam e
node attainsthrough both *A*and *B*path s inDataGuide2. T oavoidtheconfusionasin
DataGuide2, a classof DataGuidesisdefined and named *Strong DataGuide*. Themain
aspectofstrongDataGuidesisthat *eachsetoflabelpaths thatshare thesame (singleton)
targetset in the DataGuides is thesetoflabelpaths that sharesthesametargetset* inthe
datagraph.

# REFERENCES

[1]ExtensibleMarkupLanguage(XML)1.0(SecondEdition).
http://www.w3.org/TR/2000/REC-xml-20001006,WorldWideConsortium(W3C)
WorkingDraft.

[2]Goldman,R.,McHugh,J.,&Widom,J.FromSemistructuredDatatoXML:
MigratingtheLoreData   ModelandQueryLanguage,In   *ProceedingsoftheWorkshopon
theWebandDatabases*   ,June1999.

[3]XMLSchema.http://www.w3.org/XML/Schema,WorldWideConsortium(W3C).

[4]Radiya,A.  &Dixit,V.ThebasicsofusingXMLSchematodefineelements:Get
startedusingXMLSchemainsteadofDTDsfordefiningthestructureofXML
documents.http://www -106.ibm.com/developerworks/library/xml-schema/,August2000.

[5]Lassila,O.,Swick,R.R.,ResourceDescriptionFramework(RDF)ModelandSyntax
Specification, http://www.w3.org/TR/REC-rdf-syntax/,February1999.

[6]Berners -Lee,T.,WhyRDFmodelisdifferentfromtheXMLmodel,
http://www.w3.org/DesignIssues/RDF-XML.html,September1998.

[7]XMLLinkingLanguage(XLink),Version1.0,W3Crecommendation,
http://www.w3.org/TR/xlink/,20December2000.

[8]XMLPointerLanguage(XPointer),W3CWorkingDraft,
http://www.w3.org/TR/xptr,8January2001.

[9]XMLPathLanguage(XPath),Version1.0,W3Crecommendation,
http://www.w3.org/TR/xpath,16November1999.

[10] Bourret,R.,  XMLandDatabases,  http://www.rpbourret.com/xml/,February2002.

[11]Banerjee,S.,Krishnamurthy,V.,Krishnaprasad,M.,andMurthy,R.,Oracle8i       –The
XMLEnabledDataManagementSystem,    In *InternationalConferenceonData
Engineering*,F ebruary2000.

[12]Cheng,J.&Xu,J.,XMLandDB2,In       *InternationalConferenceonData
Engineering*,February2000.

[13]Rys,M.,State   -of-the-ArtXMLSupportinRDBMS:MicrosoftSQLServer'sXML
Features,BulletinoftheTechnicalCommitteeonDataEng        ineering,June2001  .

[14]Deutsch,A.,Fernandez,M.,&      Suciu,D.,StoringSemistructuredDatawith
STORED,InSIGMODConference,    June 1999.

[15]Shanmugasundaram,J.,Tufte,K.,He,G.,Zhang,C.,DeWitt,D.,&Naughton,J.,
RelationalDatabasesforQue   ryingXMLDocuments:LimitationsandOpportunities,In
*InternationalConferenceonVeryLargeDataBases*        ,September1999.

[16]Papakonstantinou,Y.,Garcia   -Molina,H.,&Widom,J.,ObjectExchangeAcross
HeterogeneousInformationSources  ,In *InternationalConferenceonVeryLargeData
Bases*,September1997.

[17]Abiteboul,S.,Buneman,P.,&Suciu,D.,DataontheWeb           :FromRelationsto
SemistructuredDataandXML,MorganKaufmann,2000.

[18]Goldman,R.&    Widom,J.,DataGuides:EnablingQueryFormulatio    nand
OptimizationinSemistructuredDatabases,In    *InternationalConferenceonVeryLarge
DataBases* ,September1997.

[19]DAML+OILWebOntologyLanguage,SubmissionrequesttoW3C,
http://www.w3.org/Submission/2001/12/,December2001.

[20]Horrocks,I.,  DAML+OIL:aReason   -ableWebOntologyLanguage,In    *Proceedings
ofEDBT* ,March2002.