

How to Develop Stylesheet for XML to XSL-FO Transformation

July, 2005
2nd Edition



Antenna House, Inc.
Copyright © 2001-2007 Antenna House, Inc.

Table of Contents

Preface	1
Step for XSL-FO Transformation	2
SimpleDoc Organization	3
Hello! World	5
Simple Example of Transforming SimpleDoc into XSL-FO	5
Stylesheet Structure	6
Block Element and Inline Element	6
XSL-FO Tree Structure	7
Developing a Practical Stylesheet	9
Printing Form Specification	9
XSL Stylesheet Organization	10
Page Layout specification	12
Page layout of Cover/Teble of contents	12
Body - Change Page Layout on Right and Left pages	14
Index - Two-Column Layout	15
Output Control as a Whole	16
Cover	17
Table of Contents	20
Templates for Creating a table of contents	20
Templates for Creating Lines of TOC	21
Counting the Nest Level	22
Setting properties according to the nest level	23
Getting page numbers	23
fo:leader	24
Example of the generated contents	25
Body	26
Templates for Processing a Body	26
Page Number Setting	27
Page Footer / Page Header	28
Page Footer Output	28
Page number	28
Running Footer	28
Page Header Output	29
Title	29
Thumb index	30
Head	32
Style Conditions of the Head	32
Templates for Processing Heads	33
Example of a Generated Title	35
Processing Inline Elements	36

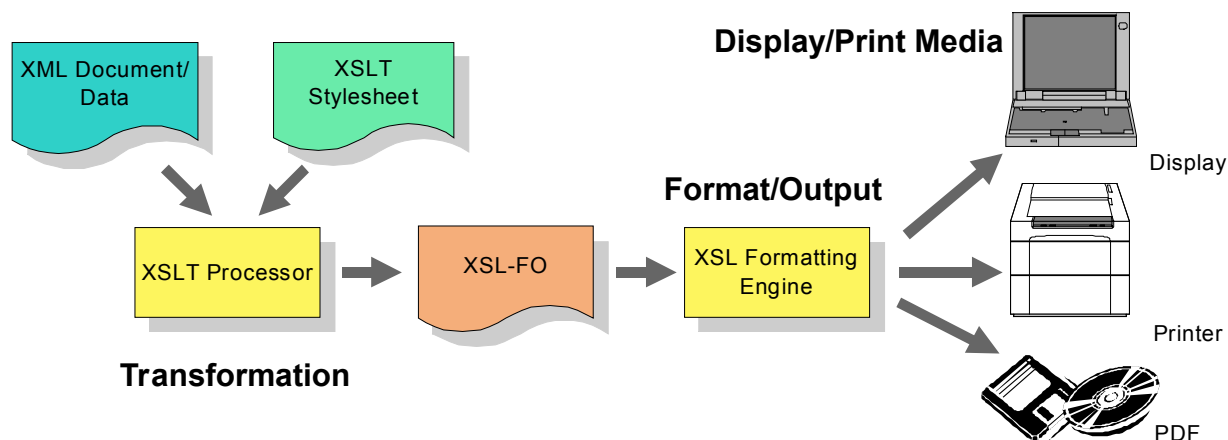
Templates that Process b, i, em, code Elements	36
a (anchor) Element	37
note Element	37
br Element	39
span Element	39
Processing Block Elements	40
p Element	40
figure Element	41
program Element	42
div Element	43
Processing Table Elements	45
Comparing SimpleDoc Table with XSL Table	45
Templates that process tables	46
Example of Table Construction	50
Processing List Elements	51
Comparing the List of SimpleDoc with the List of XSL	51
Templates that Process the Ordered List	52
Specifying the positions of label and body	53
Label Format	54
Example of Ordered List	55
Templates that Process the Unordered List	55
Specify characters for label of line	57
An example of unordered list	57
Templates that Process the Definition List	58
Templates of definition list	59
Example of a definition list	62
Functions for creating PDF	65
PDF document information	65
Creating a bookmark	65
Link setting	66
Reference to Appendix	68
Index	70
Creating keys	70
Creating an index page	70
Grouping and taking out the index elements	71
Node set output	71
Miscellaneous	73
Using modes	73
Appendix	74
Source Information	74
Update History	74
INDEX	75



Preface

Extensible Stylesheet Language (XSL) (Source Information [1]) has been brought to the attention of a wide audience as a specification recommended by W3C on 15th October, 2001 for displaying and printing the XML document. The following is the general process to transform the XML document into XSL Formatting Objects (XSL-FO) and print it.

1. Develop the stylesheet that conforms to the DTD of the source XML document to create the target output.
2. Input the XML document and the XSL stylesheet to the XSLT processor to create XSL-FO.
3. Get the target outputs such as paper-based output, PDF output, by the XSL-FO processor.



Generating XSL-FO and Display/Print by XSL Formatter

The knowledge about XSLT and XSL is necessary to develop XSLT stylesheet. The knowledge about XSLT and XSL is necessary to develop XSLT stylesheet. As for XSLT, many reference books has been published other than the specification (Source Information [2]). Probably, XSLT is already much more familiar one, since it is often used for the conversion from XML to HTML. The XSL specification (Source Information [1]) has a huge amount of contents, it has over 400 pages. It is pretty hard to understand this specification. But basically it is intended for implementers. It is not necessary for XSLT stylesheet designers to understand everything. You can fully write stylesheet by knowing some regular contents and patterns.

This report explains how to edit stylesheet which is used for transforming XML documents into XSL-FO according to the example of SimpleDoc. SimpleDoc is made for this report as a format to write a simple document. This is based on PureSmartDoc (Source Information [4]) presented by Tomoharu Asami. To make it a simple, the number of the elements is reduced and the useful functions for writing and formatting documents are added.

This report explains how to edit a stylesheet which is used for transforming a SimpleDoc document into XSL-FO. This report itself is the instance XML document of SimpleDoc.dtd, and it is then ready to be formatted by Antenna House XSL Formatter with the stylesheet which transforms the SimpleDoc document explained here into XSL-FO.



Step for XSL-FO Transformation

Now, what steps are necessary to develop XSL stylesheet? These steps are explained briefly as below.

Steps	Contents
Know the structure of the XML document	First, the information about the structure of XML source documents is required. XSLT processor can transform XML document into XSL-FO without a DTD. But the information described in the DTD such as types of elements, contents of elements, appearing order of elements and values of properties are necessary for developing a stylesheet.
Specify a printing form	This is the printing form as a final output, in other words the output specification. XSL is a formatting specification. Printing forms has various range of specifications such as sizes and layouts of printing paper, layouts of head and body, deciding whether or not to output index, table of contents, and so on.
Apply a printing form to formatting objects	After determining the specification of printing, you have to know what XSL formatting objects and properties are applied in order to print in this style. It is better to practice how to specify by referring to a simple stylesheet.
Develop an XSL stylesheet	Put the instructions to the stylesheet in order to transform XML source documents into the target printing form. Map the XML source document to XSL formatting objects that can generate the output specification. The stylesheet have the similar aspect as the general programming languages, while it may be difficult if you do not understand the feature of the XSL. ⁽¹⁾

⁽¹⁾ Refer to the definition list template in this report. In XSL, Structure for control of the conditional branch can be made, but it is impossible to assign a value to a variable. The technique to realize by calling loops recursively is necessary.



SimpleDoc Organization

The following table shows the structure of SimpleDoc treated in this report. For more detail, refer to SimpleDoc.dtd.

Element	Meaning	Definition
a group of block elements	–	p figure ul ol dl table program div
a group of inline elements	–	a note span b i em code br
doc	root element	(head, body)
head	header	(date author abstract title)*
date, author, abstract, title	header elements, date, author, abstract, title	(#PCDATA a group of inline elements)*
body	body	(chapter part section a group of block elements a group of inline elements)*
part	part	(title, (chapter a group of block elements a group of inline elements)*)
chapter	chapter	(title, (section a group of block elements a group of inline elements)*)
section	section	(title, (subsection a group of block elements a group of inline elements)*)
subsection	subsection	(title, (subsubsection a group of block elements a group of inline elements)*)
subsubsection	subsubsection	(title, (a group of block elements a group of inline elements)*)
title	title	(#PCDATA a group of inline elements)*
p	paragraph	(#PCDATA a group of block elements a group of inline elements)*
figure	figure	(title?) Specify a file by the src property.
ul	unordered list	(li*) Specify a character for label of line by the type property.
ol	ordered list	(li*) Specify format of number in the label by the type property.
dl	definition list	(dt, dd)* Specify whether to format the block in horizontal way or in vertical way by the type property.
dt	definition term	(#PCDATA a group of block elements a group of inline elements)*
dd	description of details	(#PCDATA a group of block elements a group of inline elements)*
table	entire table	(title?, col*, thead?, tfoot?, tbody) Specify whether to make auto layout or fixed by the layout property. Specify the width of the entire table by the width property.
col	column format	EMPTY Specify the number of the column by the number property, the width of the column by the width property.
thead	table header	(tr*)
tfoot	table footer	(tr*)
tbody	table body	(tr*)
tr	table row	(th td)* Specify the height of the row by the height property.
th	table header	(#PCDATA a group of inline elements a group of block elements)*

Element	Meaning	Definition
		Specify the number of the columns to be expanded across, the number of the rows to be expanded vertically by the colspan and rowspan properties. The align property allows horizontal alignment to be set to left, right, or center. The valign property allows vertical alignment to be set to top, middle, bottom or baseline.
td	table data	(#PCDATA a group of inline elements a group of block elements)* Specify the number of the columns to be expanded across, the number of the rows to be expanded vertically by the colspan and rowspan properties. The align property allows horizontal alignment to be set to left, right, or center. The valign property allows vertical alignment to be set to top, middle, bottom or baseline.
program	program code	(#PCDATA title)*
div	general block element	(title, (general block element general inline element)*) The div element expands the type by the class property.
a	anchor(link)	(#PCDATA a group of inline elements)* Specify URL as the value of href property.
note	note	(#PCDATA a group of inline elements)*
b	bold typeface	(#PCDATA a group of inline elements)*
i	italic typeface	(#PCDATA a group of inline elements)*
em	emphasis	(#PCDATA a group of inline elements)*
code	program code of the in-line elements	(#PCDATA a group of inline elements)*
span	general in-line element	(#PCDATA a group of inline elements)*
br	line break	EMPTY

The features of SimpleDoc are:

- You can start writing a document from part, also start from section. It has a flexible structure so that it can map to various kinds of documents.
- The number of the block element and in-line element are reduced to a minimum amount. The div in the block element and the span in the in-line element are defined to give various extensions.
- The br element is defined so that you can break lines inside of the lists, or the cells in the table, also inside of the paragraph without ending the paragraph.
- Output format of the list and table can be specified by the attributes.
- The problem is that the content and style is mixed.



Hello! World

Simple Example of Transforming SimpleDoc into XSL-FO

Now, we show the simplest simple stylesheet that transforms SimpleDoc to XSL-FO as follows.

Source XML Document (Hello.xml)

```
<?xml version="1.0" encoding="UTF-16" ?>
<doc>
  <head>
    <title>Simple</title>
  </head>
  <body>
    <p>Hello World!</p>
    <p>This is the first<b>SimpleDoc</b></p>
  </body>
</doc>
```

XSL Stylesheet(Simple.xsl)

```
<?xml version="1.0" encoding="UTF-8" ?>
<xsl:stylesheet version="1.0" xmlns:fo="http://www.w3.org/1999/XSL/Format"
                  xmlns:xsl="http://www.w3.org/1999/XSL/Transform">
  <xsl:output method="xml" version="1.0" indent="yes"/>

  <xsl:template match="doc">
    <fo:root xmlns:fo="http://www.w3.org/1999/XSL/Format">
      <fo:layout-master-set>
        <fo:simple-page-master page-height="297mm" page-width="210mm"
          margin="5mm 25mm 5mm 25mm" master-name="PageMaster">
          <fo:region-body margin="20mm 0mm 20mm 0mm"/>
        </fo:simple-page-master>
      </fo:layout-master-set>
      <fo:page-sequence master-reference="PageMaster">
        <fo:flow flow-name="xsl-region-body" >
          <fo:block>
            <xsl:apply-templates select="body"/>
          </fo:block>
        </fo:flow>
      </fo:page-sequence>
    </fo:root>
  </xsl:template>

  <xsl:template match="body">
    <xsl:apply-templates/>
  </xsl:template>

  <xsl:template match="p">
    <fo:block>
      <xsl:apply-templates/>
    </fo:block>
  </xsl:template>

  <xsl:template match="b">
    <fo:inline font-weight="bold">
      <xsl:apply-templates/>
    </fo:inline>
  </xsl:template>
```



```
</fo:inline>
</xsl:template>

</xsl:stylesheet >
```

Generated XSL-FO

```
<?xml version="1.0" encoding="UTF-16"?>
<fo:root xmlns:fo="http://www.w3.org/1999/XSL/Format">
  <fo:layout-master-set>
    <fo:simple-page-master page-height="297mm" page-width="210mm"
      margin="5mm 25mm 5mm 25mm" master-name="PageMaster">
      <fo:region-body margin="20mm 0mm 20mm 0mm"/>
    </fo:simple-page-master>
  </fo:layout-master-set>
  <fo:page-sequence master-reference="PageMaster">
    <fo:flow flow-name="xsl-region-body">
      <fo:block>
        <fo:block>Hello World!</fo:block>
        <fo:block>This is the first
          <fo:inline font-weight="bold">SimpleDoc</fo:inline>
        </fo:block>
      </fo:block>
    </fo:flow>
  </fo:page-sequence>
</fo:root>
```

The above XSL-FO is formatted/displayed as follows.

Hello World!
This is the first **SimpleDoc**

Stylesheet Structure

The above Simple.xml and XSL-FO show the following facts:

- A stylesheet is a set of templates. The descendant of the root element, xsl:stylesheet consists of xsl:template elements. Each xsl:template is applied so that the xxx tag of the source XML document may be processed by match="xxx".
- Formatting objects and the XML source text in each template are output to result XSL-FO tree. Then, templates that match to the descendant elements are called by an instruction of xsl:apply-templates.

XSLT processor loads the source XML document, starts processing from the root node. It finds the templates that match each node, and processes them as described in the templates. The processor processes child elements recursively, continues until the processor returns to the root element and there are no more templates to be processed.

Block Element and Inline Element

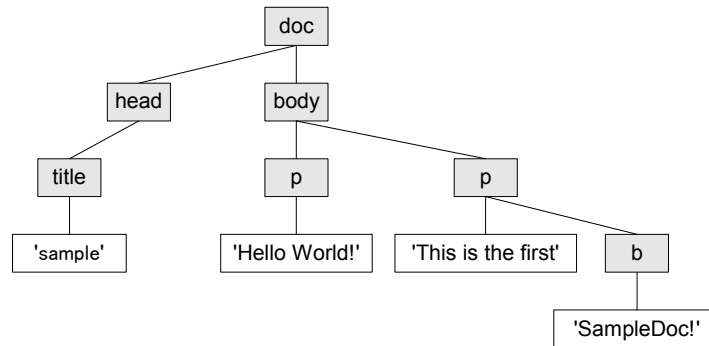
Please note how the XSL stylesheet maps block elements and inline elements in source element to formatting objects.

- In the stylesheet, p elements are transformed into fo:block objects, b elements are transformed into fo:inline objects. The base of XSL-FO transformation is to map the elements of the source XML document to either fo:block elements or fo:inline elements according to the layout instruction.
- The elements that intend to break lines by the end tag normally can be mapped to the fo:block objects. the elements of which the end tag do not intend to break lines can be mapped to the fo:inline objects. The attributes of source

elements specify properties of formatting objects. In this case, the `b` element means emphasis and property of the inline object generated from `b` is specified as bold.

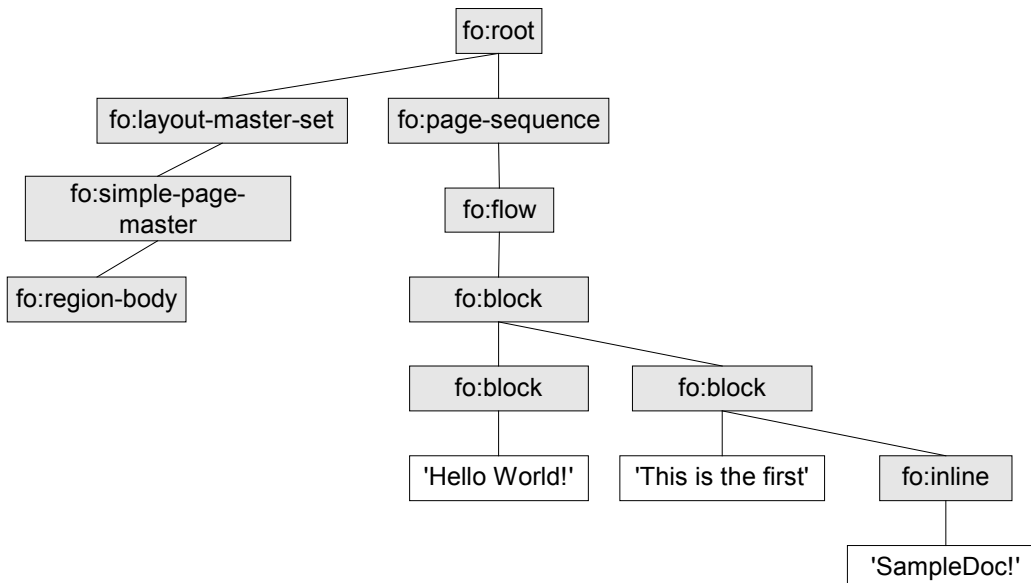
XSL-FO Tree Structure

Next, please pay attention to the XSL-FO tree structure. The following illustrates the structure of the XML document.



Hello.xml Tree Structure

In contrast, the tree structure of XSL-FO is as follows. The root of the XSL-FO tree is `fo:root`, which has two children, `fo:layout-master-set` and `fo:page-sequence`. `fo:layout-master-set` defines the page layouts and `fo:page-sequence` has a flow of contents arranged in in pages.



XSL-FO Tree after XSL Processing

`fo:layout-master-set` that defines the page layouts should precede `fo:page-sequence` (preceding-sibling) that is an actual content of pages. XSL processor processes the XML source document from the root element, seeking the templates (`xsl:template`) to be matched. Therefore, **fo:layout-master-set element should be an output from the template that processes the root element of the XML source document.** In this case, `<xsl:template match="doc">` takes this processing.

Evolution of Data Modeling for Databases

analyzed previously.

*The word Data will be used in singular throughout this article in keeping with the convention in database literature.

Since this issue of *Communications*

Basic Definitions

A **data model** is a set of concepts that can be used to describe the

structure of and operations on a database. By *structure of a database* we mean the data types, relationships, and constraints that define the “template” of that database. A data model should provide **operations** on the database that allow retrievals and updates including insertions, deletions, and modifications. Note that we will use the term “data model” to refer to the discipline for modeling data in a particular way—one that provides the building blocks or the **modeling constructs** with which the structure of a database can be described. We will use the term **application model** to refer to the description of data for a particular database. For example, the relational model is a data model. The definition of a particular database for the personnel application at company X will be called an application model of that database which uses the relational model. Application analysts often refer to the latter as a data model, causing confusion.

In any application, it is important to distinguish between the *description* of the database and the *database itself*. The description is called the **database schema**. A database schema is designed for a given set of applications and users by analyzing requirements. It is described by using a specific data model that provides the modeling constructs in the form of a language syntax or diagrammatic conventions. In some tools or DBMSs, the data model is not explicitly defined but is present implicitly in terms of the features present. Schemas usually remain relatively stable over the lifetime of a database for most applications. For dynamic environments, such as computer aided design (CAD), or computer-aided software engineering (CASE), the schema of the product being designed may itself change. This is referred to as **schema evolution**. Most DBMSs handle a very limited amount of schema evolution internally.

During the process of database design, the schema may undergo transformation from one model

into another. For example, the schema of the personnel database may be initially described using the entity-relationship data model in the form of an ER diagram. It may then be mapped into the relational data model which uses structured query language (SQL)—an emerging standard, to define the schema. The entire activity of starting from requirements and producing the definition of the final implementable schema in a DBMS is called **schema design**. The DBMS is used to store a database conforming to that schema; it allows for database **transaction processing** in which a transaction is a unit of activity against the database that includes retrieval and update operations against the database. A transaction must be performed in its entirety, leaving a “permanent” change in the database or may be aborted, performing no change at all.

The actual database reflects the state of the real world pertaining to the application or the “miniworld.” It must remain in conformity with the miniworld by reflecting the actual changes taking place. The data in the database at a particular time is called the “database instance” or the “database state.” The actual **occurrences** or **instances** of data change very frequently as opposed to the schema, which remains static.

A data model in the database parlance is associated with a variety of languages: data definition language, query language, data manipulation language, to name the important ones. The **data definition language (DDL)** allows the database administrator or database designer to define the database schema. The DBMS has a compiler to process the schema definition in DDL and to convert it into a machine-processable form. This way, a centralized definition of the application model is created, against which a number of applications can be defined. **Data manipulation language (DML)** is a language used to specify the retrieval, insertion, deletion, and modification of data. DMLs may be divided broadly into

two categories: declarative and procedural. The former allow the user to state the result (of a query) that he or she is interested in, whereas the latter require one to give a procedure for getting the result of the query. The nature of the language also depends on the data model. For example, while it is possible to have a declarative language for a model such as the relational model, the language for the network data model is “navigational,” requiring the user to state how to navigate through the database, and thus is inherently procedural. Either type of language may be used in a stand-alone interactive fashion as a **query language**. Languages for data models can also be distinguished in terms of whether or not they are record-at-a-time or set-at-a-time. Record-at-a-time processing requires an elaborate control structure typically provided by a **host programming language** within which the DML commands or verbs are embedded. Set-oriented processing regards data as sets of elements (e.g., sets of tuples in the relational model) and provides for operators that apply to these sets, generating new sets. There is now a movement toward providing languages which seamlessly integrate the capability to provide general-purpose computation and special-purpose data manipulation against the database in a single language. These are called **database programming languages (DBPLs)** [6].

Scope of Data Models

In the traditional sense, data models used for database schema design have been limited in their scope. They have been used to model the *static* properties of data including the following:

Structure of data: The structure is expressed in terms of how atomic data types are aggregated into higher-order types. Furthermore, the models express relationships among these aggregates. The early models tended to be record-oriented, the basic aggregate struc-

ture being a record type consisting of data element types or field types. The database schema consists of a number of record types that are related in different ways. The hierarchical data model organizes the record types in a tree structure, whereas the network data model organizes a database schema with the record types as the nodes of a graph. The limitations of the record-based view of data modeling are discussed in [23]. The relational model introduced a set-oriented view of data modeling [16], and currently, the object-oriented view which structures a database in terms of objects and interobject interactions is gaining popularity [3]. We discuss these approaches in greater detail later in this article.

Constraints: Constraints are additional restrictions on the occurrences of data within a database that must hold *at all times*. Data model constraints serve two primary goals:

- Integrity: Integrity constraints are the rules that constrain the valid states of a database. They arise either as properties of data, or as user-defined rules that reflect the meaning of data.
- Security and protection: This applies to restrictions and authorization limitations that are applied to a database to protect it from misuse and unauthorized usage.

Constraints can be visualized at different levels:

- a) **inherent constraints** pertain to the constraints that are built into the rules of the data model itself. For example, in the entity-relationship model a relationship must have at least two participating entity types (depending on the variant of the model used, the same entity type may be used twice).
- b) **implicit constraints** are constraints that can be specified using the DDL of a data model to describe additional properties. They are expected to be automatically enforced. An example is a mandatory participation constraint in the

entity-relationship model, which states that a specific entity must participate in a particular relationship.

c) **explicit constraints** are application-dependent constraints that specify the semantic constraints related to an application. These are the most general and difficult constraints to specify in full detail. There is a general trend to capture as much "application behavior" information within a database as possible in order to capture it in a central place. The 4GL movement is aimed at capturing this application semantics at a high level. Today's DBMSs are not equipped to handle such constraints easily, however.

Another dimension of constraint modeling is to capture state transition rather than just static state information. This gives rise to **dynamic constraints** which are stated in terms of what types of changes are valid on a database. An example is: "the salary of an employee can only increase." Both static and dynamic constraints allow us to define whether a database is consistent. Whereas static constraints can be evaluated on a "snapshot" of a database to determine its validity, the dynamic constraints are much more difficult to enforce, since they involve blocking/preventing a state transition at "run-time."

Other parameters of data models: A data model for a database may also specify some additional details relevant to the use of the data. One possible feature is "**distribution parameters**." These relate to the fragmentation of data in terms of how data is stored as fragments. In the relational model it is customary to refer to "horizontal" and "vertical" fragments. The former contain subsets of the data occurrences of a relation that meet some predicate condition, while the latter refer to a subset of the attributes of data for the whole relation. In a relation called ORDERS, each horizontal fragment may contain orders that are shipped from one warehouse, whereas ORDER may be vertically

fragmented into shipping information and billing information. **Security** is another feature that may be built into a data model at different levels of granularity. Yet another feature is **redundancy**, which is hard to model explicitly; it may be captured in the form of specification of explicit copies or overlapping data. A model must allow for specification of features such as **keys** which uniquely identify data; it may also have a way to specify physical parameters such as **clustering** or **indexing** (e.g., B+ tree index on a certain field) as a part of the application model specification.

Views: In database modeling a **view** is a perceived application model as defined by a user or an application group. A view is another mechanism by which an application can record its specific requirements, in addition to the explicit constraints mentioned previously. Most data models provide a language to define views: it may be coincident with the query language, whereby the result of a query is defined as the view. Typically, a view is constructed when there is a request to retrieve from it, rather than "materializing" the view because the latter creates redundant data.

Toward Joint Data and Application Modeling

To give adequate support for the modeling of dynamic application environments, another school of thinking combines the functional analysis, which is typically the focus of software engineering during information system design, with conventional data modeling [11, 26], giving rise to a joint analysis and modeling of application and data requirements. We have also advocated this view during conceptual database modeling in [8]. In terms of the present issue of *Communications*, the preceding approach seems most relevant and significant. A top-down structured design methodology for both data and process modeling using extended ER and data flow diagrams respectively is