

Structured Document Storage and Refined Declarative and Navigational Access Mechanisms in HyperStorM^{*†}

Klemens Böhm, Karl Aberer, Erich J. Neuhold, and Xiaoya Yang
GMD-IPSI, Dolivostraße 15, 64293 Darmstadt, Germany
{kboehm, aberer, neuhold, yangx}@darmstadt.gmd.de

Abstract

The combination of SGML and database technology allows to refine both declarative and navigational access mechanisms for structured document collection: with regard to declarative access, the user can formulate complex information needs without knowing a query language, the respective document-type definition or the underlying modelling. Navigational access is eased by hyperlink-rendition mechanisms going beyond plain link-integrity checking. With our approach, the database-internal representation of documents is configurable. It allows for an efficient implementation of operations, because DTD knowledge is not needed for document structure recognition. We show how the number of method invocations and the cost of parsing can be significantly reduced.

Keywords: OODBMSs, SGML, document query languages, navigation

^{*}‘HyperStorM’ is an acronym for ‘Hypermedia Document Storage and Modelling’

[†]This article has been accepted for publication in VLDB Journal, Volume 6 (1997),

1 Introduction

1.1 Objective of this Work

With open systems, such as the World-Wide Web (WWW), and document-exchange formats where markup can be placed within the document at ease, notably HTML [HTM], there neither is control over the structure of individual documents nor over consistency of the document collection. With the combination of SGML ('Standard Generalized Markup Language') [ISO86] and database technology, on the other hand, one can specify the logical structure of documents, make assumptions about their structure and ensure the consistency of the document collection. This allows to refine both declarative and navigational access mechanisms.

Declarative access. With regard to declarative access to a document collection, the prevailing question in literature in the past few years has been: which is the most expressive query language? However, the more expressive the query language, the more complex it is. In addition to having an interface allowing the formulation of complex queries, ease of use of the search mechanisms has a high priority for a lot of users. In this article, we describe mechanisms for declarative access, so-called *query templates*. They allow the user to search with a fairly high expressive power without having to know a query language, the respective document-type definition, or the underlying modelling. Different query templates can be made available for individual user groups.

A query template is a document-type-specific query form that is generated automatically from the corresponding document-type definition and an additional query-template specification. With a query template, only a limited set of element- and attribute types is made available for querying. Hence, the query template can be seen as a view mechanism. The query-template specification is administered by the database application. This

ensures consistency of the query template with the respective document-type definition. In consequence, only meaningful queries can be generated from the user input.

With regard to query templates, we exploit the fact that there is a document-type definition, i.e., a schema that documents of a certain type must conform to. On the contrary, if HTML-documents were to be administered, the following problems would be in the way of having query templates and making proper use of them:

- **No support for individual document types.** While trying to provide coverage for all kinds of documents, HTML is highly generic. SGML, however, has been developed to reflect the particularities of different application scenarios by using different document types. Queries can be much more specific if the document type may be taken into account.
- **Meta-information is not available or may not be consistent.** With appropriate document-type definitions, any kind of metadata¹ might be seamlessly included into documents to cover users' information needs [KSS95]. With HTML 3.2, consistency of metadata in different documents is not ensured. The HTML-DTD does not contain any guidelines how meta-information should be modelled.
- **Difficult to consider documents' internal structure.** As opposed to SGML-DTDs specified for a particular application scenario, HTML allows to model documents' logical structure only in a generic way. In consequence, search mechanisms are only a subset of the ones for full-fledged SGML documents.

¹Throughout this article, the terms 'metadata' and 'metainformation' are used interchangeably.

Navigational access. Another aspect of applying both SGML and database technology is support for navigation. The following problems are addressed:

- **“Navigation Failure”.** Navigation to a document may fail for a variety of reasons: links to other documents may be outdated, i.e., the document referenced may have been removed or moved to another location in the meantime. Navigation to a document referenced may not be possible because the reader is not entitled to access it, or because the client is unable to present the (multimedia) document. E.g., a viewer for the respective format may not be available at the client site.
- **“Expectations vs. Actual Content”.** Frequently, the reader has no meaningful information on the document referenced by a hyperlink within the current document or on the set of documents to be reached by following a hyperlink anchor. Experience shows that, quite often, expectations awakened by the context of a link anchor are not matched when actually viewing the document. In other words, in order to find out more about the target of the link, there seems to be no alternative to bringing the document to the client site.
- **“Lost in Hyperspace”.** Traversing several hyperlinks may have the effect that the reader loses orientation in the document collection. Being “lost in hyperspace” [Con87] stems from the fact that hypermedia objects, e.g., the WWW pages, are arranged in a non-linear way. The reader may have to inspect several such objects to satisfy his information need. He may not be able to find all the relevant objects by navigation or to avoid the irrelevant ones.

These problems are alleviated by our approach. With WWW-based access to the document base, a document is converted to HTML according to a conversion specification before presentation. This specification is contained

alliances, goose and duck pate de foie gras, truffles...the [Perigord](#) (LANGUAGE: FRENCH) has "a fragrant soul" Sarlat with its medieval authenticity, is its capital.

Figure 1: example for hyperlink rendition

in a style-sheet. Conversion takes place within the database. Namely, rendition of a hyperlink anchor pointing to another document may depend on that document's characteristics.² In particular, metadata on the document referenced as well as documents indirectly referenced may be used for link anchor rendition. As a simple example, consider a document collection where the language of the document is made explicit in the document header, i.e., a document with title 'Perigord' contains '<LANGUAGE>FRENCH</LANGUAGE>'. If document 'Perigord' is referenced by another document, the respective hyperlink anchor in the layouted version of the document may look as in Figure 1. This alleviates the 'Expectations-vs.-Actual-Content' problem. In this example, if the reader does not speak French, he knows that he does not want to access the document.

With hyperlink rendition, information on the document directly referenced, but also on documents indirectly referenced, may be taken into account. Information to be displayed within the hyperlink is specified in the database-query language. To speed up the process of retrieving a document from the database, HTML conversion results may be materialized in the database.

If we again compare the administration of HTML documents to our approach, we face similar problems as with declarative access. Meta-information may not be available for document rendition, or one cannot tell which meta-information has been made explicit within documents. Inconsistent repre-

²Throughout this article, a document containing the anchor of a link referencing another document will be called *referencing document*, the other document will be referred to as the *document referenced*.

sensation of metadata likewise is a problem.

1.2 Database-Internal Representation of Documents

An efficient database-internal representation of structured documents is a prerequisite to realize the functionality that has been outlined so far. We advocate a hybrid database-internal representation of documents. Only “big” elements are represented by individual database objects. Different “small” elements, so-called *flat elements*, are mapped to the same database object. The structured representation of documents is advantageous to allow fine-grained modifications of documents in the database and to reflect the semantics of hypermedia document components, whereas performance of certain basic operations, such as insertion of a document into the database, is better with an unstructured representation [BAK95]. With regard to the mapping of flat elements to database objects, there are design alternatives that are described in this article.

Flat elements are completely marked up in the database. This allows a more efficient implementation of methods operating on the document structure, as compared to parsing the document fragment (see Section 4). Another aspect is that elements can be disambiguously identified by means of a logical object identifier. If, finally, methods are part of the database, as is the case with an OODBMS, method invocations are relatively expensive.³ It is advantageous to take this into account when implementing the methods.

The database-internal representation of documents is configurable. The initial configuration is described by means of an SGML document, an instance of the so-called *super-DTD*. The configuration may be partially al-

³Our terminology is different from the one used in [Sto96]. Object-oriented DBMSs in our terminology essentially correspond to object-relational ones in [Sto96]. There is no counterpart in our terminology for object-oriented DBMSs in that terminology, which basically are persistent object-oriented programming languages without sophisticated declarative access mechanisms.

tered at runtime. Using SGML to describe the startup configuration has the following advantages:

- The consistency of this specification is checked by means of an SGML parser.
- To specify the initial configuration, the specification document is inserted into the database in the same way normal documents are. No differentiation is necessary from the user's perspective.
- The DTD-designer can specify the startup configuration using a mechanism he is familiar with.
- The DTD-designer may use SGML tools to fulfill his task.
- The mechanisms described in this article for querying the document collection, notably query templates, are also available for document types as documents. In other words, the document type is explicitly available, in analogy to database-schema information contained in a data dictionary.

The remainder of this article has the following structure: in the following section, related work is reviewed. Section 3 contains a brief review of SGML concepts and introduces some notions that are relevant in this context. In Section 4, we describe the database-internal representation of documents in our framework. Section 5 describes how declarative access functionality can be extended by using SGML and database technology, Section 6 focuses on navigational access. Section 7 concludes the paper.

2 Related Work

We have classified related work as follows: related work with regard to document modelling, related work with regard to declarative access, and related work with regard to navigation and HTML conversion.

Related work with regard to document modelling. In order to represent the documents' logical structure using an OODBMS, it seems feasible to carry out a 1:1 mapping from elements to database objects. Furthermore, there would be database classes, so-called *element-type classes*, corresponding to an element type from the DTD. However, if this logical view is identical with the physical representation, the following problem will arise: the duration of certain basic operations such as inserting documents into the database or retrieving documents from the database is almost directly proportional to the number of database objects that are created or retrieved, respectively. This may not be acceptable, as others have observed, too [NBY95]. - An alternative seems to be the approach described in [ACM93, ACM95]. They consider structured data whose physical representation is flat, in particular data within files. If the structure is needed for, e.g., query evaluation or updates, the document is parsed, and objects in main memory are generated. Our work differs from theirs in the following respects:

- With our approach, a document is not necessarily represented by one file. Rather, the document may be physically fragmented in the database, and the fragments' logical structure can be recognized using, e.g., the techniques described in [ACM93, ACM95]. The database-internal representation is configurable. One advantage is that a finer granularity is possible with regard to concurrent write access to the document, using concurrency control mechanisms provided by the DBMS.
- With our database-internal representation, the DTD is not needed to recognize documents' logical structure. More specific techniques than parsing the document are applied. The advantages will be pointed out in Section 4.

- An element has an OID whose lifetime is independent from the existence of a corresponding object in main memory.

Related work with regard to declarative access. The following issues are of interest with regard to work on declarative access to document collections having originated in the database community [Mac91, C⁺94, MMM96, QRS⁺95, B⁺94, ST94, O⁺95].

- With our approach, expressiveness of the query language is achieved by using methods of the database schema, together with OQL query mechanisms [Cat94]. Compared to other approaches, the expressive power is higher, while, on the other hand, it is not necessary to extend the query algebra. Our approach allows for full-fledged information-retrieval functionality (IR-functionality) [VAB96], which is different from search on a syntactic level [SM83], as well as search on documents' physical characteristics.
- With many of the above references, information with regard to documents' database-internal representation is incomplete or missing.
- To the best of our knowledge, work cited above does not contain any counterparts to query templates or the hyperlink rendition mechanisms described in this article.

In more detail, work described in [C⁺94] is based on OODBMS-technology. They have extended the underlying query algebra to reflect notions such as the lengths of paths. By using OQL that allows inclusion of methods into query statements, and having an adequate set of methods as part of the database schema, the same expressiveness can be obtained without extending the query algebra. In [B⁺94], a datatype 'structured text' is introduced to be integrated into relational database systems and an extension of SQL is defined. To facilitate updates, the approach is to map SGML structures to

tables, but conformance to the DTD remains to be ensured. The PAT-query algebra [ST94] lacks certain features such as the notion of position, querying according to documents' secondary structure, and aggregation. Further, only elements can be retrieved. It is, however, independent of the data-model and will be dealt with in Section 5 again. While in [O⁺95] a user interface for an SGML-/HyTime-document database has been realized, work seems to have been centered around one particular document type.

In [YA94], a coupling of a DBMS and a text search engine is described. There, documents' internal structure is not modelled within the database. The text engine used there does not support the notion of vagueness. The need for IR-functionality, e.g., ranking, is acknowledged in [SDAMK95]. Their objective is to build an integrated system providing both database- and IR functionality. Details about documents' internal representation are not revealed. We for our part have realized a loose coupling between the OODBMS VODAK [VML95] and the IR-system INQUERY [CCH92] to make IR functionality available for database content. With a loose coupling, we will be able to rather easily incorporate improved IR-functionality whenever it becomes available.

An objective of others, e.g., [QRS⁺95, MMM96], is to provide declarative access mechanisms for open-ended systems where assumptions about the data's structure cannot be made, notably the WWW. Even though WWW-related issues currently draw a lot of attention, the question how to exploit consistency of the document collection in controlled environments remains to be relevant.

Related work with regard to navigation and HTML conversion.

With Hyper-G [AKM95], a principal objective is to ensure hyperlink consistency. The idea is that there is a link database. It contains the information which hyperlinks exist between documents. The advantage, as compared to

the current status of our work, is that there is no confinement to the content of one database. This is reached by giving up some of the individual information servers' autonomy. Hyper-G is not modular, but, rather, can be seen as "another web", as, for example, proprietary browsers have to be used. In our context, a mere link database would not be sufficient, as arbitrary information on documents referenced can be requested. An open question is whether people are willing to take into account the additional overhead of "a WWW without dangling references". On the other hand, the need to ensure consistency of local document collections clearly exists [S⁺94]. - Conversion of SGML documents to HTML is the topic of [Fre]. The notion of *location grammar* is introduced as a means to specify context-sensitive transformation of element types. It seems that, there, context-sensitivity refers to documents' hierarchical structure, but not to other documents. The topic of [kbo94], similarly, is structured document handling in the internet. They argue that it is the SGML document that should be delivered to the client to facilitate so-called *document post-processing*. Trivially, our database server can also return the original SGML documents. If the DTD allows for it, hyperlink rendition, as outlined above, is still feasible. Annotation servers contain information on WWW documents that may be provided by others. Instead of directly bringing the document to the client, the document goes through the annotation server, and relevant information is added [RMW94].

DSSSL [ISO96] is an expressive language to specify document transformation. The standard specifies a structured representation of documents; conversion is based on that representation. The standard does not deal with the question how to efficiently carry out such a conversion if documents are within a storage system, and if characteristics of documents referenced are taken into account. With our approach, characteristics of the documents referenced can be reflected. Rendition mechanisms for hyperlink anchors likewise are expressive, and are identical with our declarative access mech-

```

<agenda author=Aberer>
  <header><language>English</language>
    <subject>future research topics for the department</subject>
    <location>...<date>...
    <invited>
      <name>Fischer</name><name>Chen</name>...
    </invited></header>
  <programme>
    <item>brief review of present funding situation:  in 1996/97 ...</item>
    <item>problems with diploma thesis students:  due to the fact ...</item>...
  </programme>
</agenda>

```

Figure 2: sample SGML document of type ‘Agenda’

anisms.

3 Modelling Meta-information with SGML

The practical relevance of SGML has considerably increased in the recent past. This is possibly due to the close connection between SGML and HTML, the format of WWW documents. Within SGML documents, the logical document components, so-called *elements*, are made explicit by means of markup. The document fragment from Figure 2 is an example of a marked-up SGML document. ‘<item>’, ‘</item>’ identify (the start position/the end position of) an element of type *item*. ‘<’ is the *start tag open (STAGO)*, using SGML terminology, ‘</’ is the *end tag open (ETAGO)*, and ‘>’ is the *tag close (TAGC)*. It is an important aspect of SGML that markup may not arbitrarily be chosen and be placed within documents. Rather, for each document type, a so-called document-type definition (DTD) has to be provided. It specifies which element types

may occur in a document, and how elements may be arranged within a document. A DTD is a grammar. The **agenda**-DTD is contained in Figure 3. Examples of element types from this DTD are **header**, **programme**, and **item**. The regular expression specifying the admissible content of an element of the respective type is referred to as *content model*. For instance, ‘(header, programme)’ is the content model of **agenda**: an **agenda**-element contains a **header**-element, followed by a **programme**-element. The expression ‘+(keyword)’ is an example of an *inclusion model*. It specifies that the structure within the brackets, in this case an element of type **keyword**, may occur arbitrarily within an element of type **programme**. For instance, the element ‘<programme><item>brief review of present <keyword>funding</keyword> situation ...</item> ...</programme>’ conforms to the DTD from Figure 3. *Exclusion models* are also available to forbid such inclusions in a subtree of the subdocument. If the element-type definition of **item** was

```
<!ELEMENT item (#PCDATA) -(keyword)>
```

the above sample element of type **programme** would not conform to the sample DTD any more. - CDATA and (#PCDATA) are terminal element types comparable to the data type STRING. - Elements may be furnished with attributes. Again, the attributes cannot freely be chosen, but must be contained in the DTD. #IMPLIED means that a value does not have to be assigned to the respective attribute.

Classifying element types. Element types can be categorized by the role of their instances within the documents [HHM94]. This classification is important as access patterns, notably with regard to declarative access, are different for the individual categories.

- **Structural element types.** Markup of such elements is used to identify documents’ logical structure. Examples from the sample DTD

```

<!DOCTYPE agenda [
<!ELEMENT agenda      (header, programme)>
<!ATTLIST agenda      author CDATA #IMPLIED>
<!ELEMENT header      (language?, (location|roomno), date, invited)>
<!ELEMENT programme    (item)+      +(keyword)>
<!ELEMENT (language|location|roomno|date|name|keyword)    (#PCDATA)>
<!ELEMENT invited      (name)+>
<!ELEMENT item          (#PCDATA)>
]>

```

Figure 3: Sample DTD (Document Type ‘Agenda’)

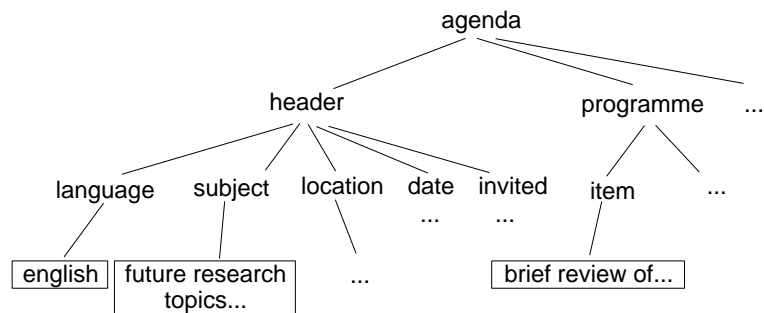


Figure 4: sample document’s hierarchical structure

are `programme` or `item`.

- **Non-structural element types.** Non-structural elements are individual words or short sequences of words within structural elements' content having a particular role, e.g., element type `keyword` in the `agenda-DTD`. In other words, markup of non-structural element types is used to make explicit the meaning of words within text. In most cases, non-structural elements are not bound to structural element types, but may occur rather freely within the text.
- **Informational element types.** Informational elements are meta-information. While non-structural elements occur within actual document text, informational elements do not occur within structural elements' textual content. Rather, they tend to be contained in a document header. Typically, informational element types that do not have an internal structure could also be modelled as SGML attributes, while non-structural element types cannot. An example of an informational element type is element type `language` from the DTD in Figure 3.

Using this categorization of element types according to their roles, we are now in the position to describe how meta-information can be modelled with SGML.

- Informational elements are metadata.
- The markup of structural elements and non-structural elements is meta-information. This is different from informational element types where the elements themselves are meta-information.
- Elements, normally structural elements, can be furnished with attributes, as described above. The attribute values are meta-information.

Furthermore, the document-type definition itself can be seen as metadata. Namely, the different ways to represent metadata, as described above,

must be complemented with the type definition for meaningful interpretation.

Further SGML mechanisms. The SGML concepts that have been described in this section are merely a subset of SGML. It is the subset for which support is described in the following. With our approach, SGML entities and marked sections are resolved by the parser and do not occur any more within the document in the database. Hence, the mapping of a document to the corresponding database content is not loss-free. Furthermore, notations and the SGML link mechanism are not supported.

4 The HyperStorM Database Application Framework to Administer Structured Documents

The structure of this section is as follows: the database-internal representation of documents is described in the following subsection; configurability mechanisms are described in Subsection 4.2. In the last subsection, the transformation algorithm from documents to their database-internal representation is presented.

4.1 Reflecting the SGML Information Model

This subsection covers design decisions and issues with regard to the database-internal representation of structured documents within the database.

1. Hybrid database-internal representation for documents: some elements are represented by individual database objects, while others, the flat ones, are not. This representation is subject to configuration for the particular document type, and the respective configuration mechanisms will be described in the sequel.
2. Flat elements are completely marked up within the database.

3. Elements have a logical OID whose life cycle is independent of the existence of corresponding (C++-)objects in main memory.
4. The query language of our system is OQL, together with methods from the database schema. With methods as part of the query language, expressiveness of the declarative access mechanisms is naturally higher than with other approaches. Method invocations are costly with methods being part of the database. This must be reflected with their implementation.

To reflect documents' internal structure, not only documents, but also document components are explicit within the database. A differentiation between flat and non-flat elements is made (cf. Section 1.2). Database objects corresponding to flat elements are *flat objects*. The string representation of a flat object's elements is the *flat string of the database object*. As an example of documents' database-internal representation, consider the document from Figure 2. One out of many representations that are possible within the database is given in Figure 5. With that particular configuration, **language** and **subject** are examples of flat element types. '`<language>English...</invited>`' and '`<programme><item>brief...</programme>`' are examples of flat strings. '`<language>English... </invited>`' is the flat string of the language- and the subject-element. The hybrid database-internal representation facilitates modifications of document fragments and better reflects the semantics of hypermedia document components. It reduces the negative impact of a structured physical representation with regard to performance [BAK95].

Structure recognition of flat elements. As just explained, elements in the database can either be flat or non-flat. While, with SGML, it is allowed to omit markup if the document structure can be disambiguously recognized by means of the DTD, document fragments within the database

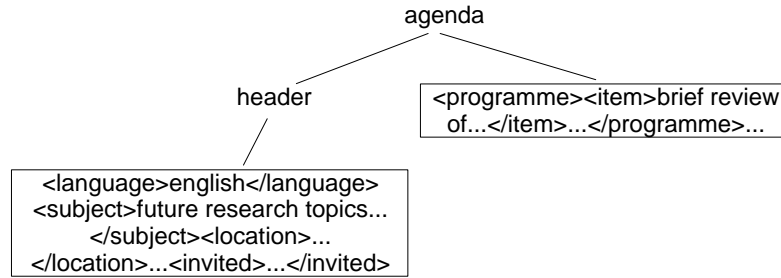


Figure 5: possible physical representation of ‘Agenda’-document

are completely marked up. Markup that may have been omitted from the original documents is added. Consequently, the document structure can be recognized without the DTD (see the top right fragment in Figure 6 as an example, as opposed to the top left one). Simple linear access operations are sufficient. The advantages of not using the DTD are the following ones:

- If the DTD was used for structure recognition, it might seem feasible to construct a fragment-DTD on-the-fly. However, DTDs are not context-free due to inclusions and exclusions. Hence, to construct the fragment DTD, one would have to inspect the inclusion and exclusion models of the elements the current flat elements are contained in. This requires a number of access operations to database objects that are unnecessary with our approach. As an example, consider the following clipping from a DTD.

```
<!ELEMENT A (C)* +(G)>
```

```
<!ELEMENT B (C)*>
```

```
<!ELEMENT C (D?, E+)>
```

In order to construct a fragment-DTD for an element *c* of type *C*, in particular, the inclusion model of type *C*, one must check if *c* is contained in an element of type *A* or *B*. In the first case, the fragment-DTD must reflect that *c* may contain an element of type *G*, as opposed

to the second case.

- Structure recognition is more efficient without the DTD: if documents of different types are in the database, it is not necessary to look up the respective type first.
- With a DTD-based approach, a flexible fragmentation of documents in the database is not possible in practice. As an example, consider the bottom left database object in Figure 5. The corresponding document fragment ‘<language>...</invited>’ does not have a root element. In consequence, either a document fragment suitable for parsing would have to be constructed first. ‘<header><language>...</invited> </header>’ would be such a parseable fragment. However, this requires access to at least one more database object, namely the header-object. If this object contained non-flat elements, further database-access operations would be necessary. Alternatively, concatenation of flat elements to build a flat object would have to be forbidden, i.e., the **language**-element, the **subject**-element etc. would be separate database objects. But this may lead to a big increase in the number of database objects.

Object identifier. Object identity is an important notion in object-oriented modelings. The necessity of OIDs for both flat and non-flat elements introduces a logical and a physical object level. The logical view remains to be the one that there is an object corresponding to each element. On the physical level, however, this is not the case. A logical OID consists of a physical OID and the STAGO-position within the corresponding flat string, i.e., the byte offset. If the respective element is a non-flat one, the offset is -1. The DBMS has been extended so that it can transparently support method invocation on objects identified by logical OIDs. With message calls, the DBMS resolves logical OIDs and dispatches them to the physical objects.

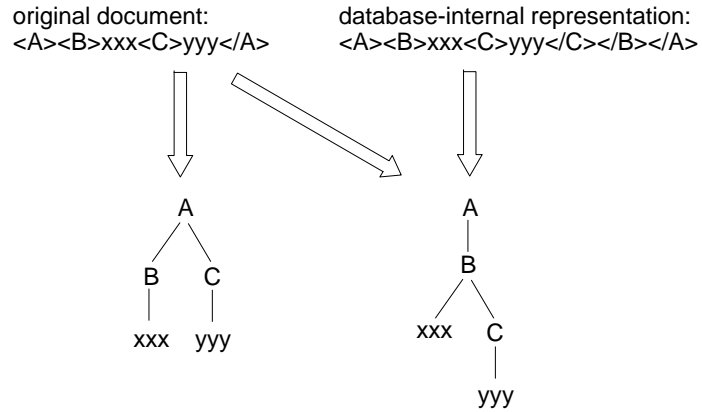


Figure 6: inferring the document structure from complete and incomplete markup

In the parsing approach described in [ACM93, ACM95], object identifiers are available only as long as the corresponding structure in main memory exists.

Classes and methods of the database schema. The following classes are part of the database schema.

ELEMENT - The physical database objects representing the document structure are instances of this class, e.g., the nodes in Figure 4.

ElementType - For each element type from a DTD, there is an instance of the class.

DTD - An instance of this class corresponds to each DTD currently supported.

Document - For each document, there is a corresponding object.

The methods for elements include the following ones:

hasTextualContentRegex (r: REGEX): BOOL

hasAttrValueRegex (attrName: STRING, r: REGEX): BOOL

```

getIRSValue (q: STRING): REAL
isContainedIn (e: logicalOID): BOOL
getReferencedElements (attrName: STRING): {logicalOID}
getAttrValue (attrName: STRING): STRING
getSize (): REAL

getAll (elementTypeName: STRING): {logicalOID}
getFirst (elementTypeName: STRING): logicalOID
getElementText (): STRING

```

Instances of `ElementType` have method `getElements(): {logicalOID}`.

Method `hasTextualContentRegex` returns `TRUE` iff `r` is contained in the target element's textual content. Method `hasAttrValueRegex` returns `TRUE` iff the value of attribute `attrName` contains `r`. `getIRSValue` returns the belief value of the element's textual content with regard to IRS-query `q`, as computed by the underlying IRS. `isContainedIn` returns `TRUE` iff the target element is contained in the parameter element `e`. If attribute `attrName` of the target element of `getReferencedElements` is of type `IDREF(S)`, the logical OIDs of the elements referenced (within the same document) are returned. Otherwise, the empty set is returned. `getAttrValue` returns the value of attribute `attrName`. `getAll` returns all elements of type `elementTypeName` that are contained in the target element; `getFirst` returns the first element of type `elementTypeName` (in preorder) that is contained in the target element. `getElementText` returns the target element's textual content. `getElements` returns the logical OIDs of all elements of the type.

With regard to methods `hasTextualContentRegex` and `getIRSValue`, some comments are appropriate. Data administered by a storage system may be subject to different paradigms. With some element types, one wants to search their instances with exact mechanisms, i.e., by means of pattern matching on the syntactic level such as regular expression search. Such search mechanisms are in place for element types such as `SURNAME` or `PART_NO`.

In this case, method `hasTextualContentRegex` should be used. On the other hand, the objective of IR search is to cover the user's information need by going beyond the syntactic level. Results of IR queries are never precise and may differ from system to system, as the content of a piece of text may be seen differently by different systems. In the IR context, it is too undifferentiated to merely say 'The document is relevant.' or 'The document is not relevant.'. Rather, relevance is expressed by means of a belief value b such that $b \in [0;1]$. The belief value is the probability that the document is relevant with regard to the query, as computed by the system. As a rule of thumb, IR mechanisms for text only work well for texts containing more than 20 to 30 words, they do not work for individual words or short sequences of words. In consequence, search on the syntax level makes sense for informational element types, i.e., metadata, while IR search mechanisms should only be applied to structural element types, i.e., raw data.

Not only individual concepts, but also complex query terms in the IRS-query language may be parameter of method `getIRSValue`. With INQUERY being the underlying IRS, parameters such as '`#and(HyTime, MHEG)`', '`#not(Java)`', or '`#uw10(SGML, HTML)`' can be processed. (The last expression specifies that 'SGML' and 'HTML' must occur within a window of 10 words.) In combination with other search mechanisms, this gives rise to a powerful search functionality.

The set of methods reflects our practical experiences and is now stable. Methods `hasTextualContentRegex`, `hasAttrValueRegex`, `getIRSValue`, `isContainedIn`, `getReferencedElements`, and `getSize` are necessary to formulate queries corresponding to terms in the extended PAT language, and none of these methods can be omitted without subtracting from expressiveness (cf. Section 5.2). Method `getAll` is an example of a method that is needed for more efficient query evaluation, as compared to `isContainedIn`. `getFirst` and `getAttrValue` go beyond the expressiveness of the extended PAT language.

Example of method implementation. Method `next` identifies the right sibling of the target element in the logical document structure. In the sample document from Figure 4, the next element of the subject-element is the location-element. Method `next` makes use of method `getPositionOfETAGO`. `getPositionOfETAGO` returns the (byte offset) position of the end tag open which corresponds to the start tag whose (byte offset) position is the method parameter. The method illustrates that, with our database-internal representation, operations on documents' logical structure are feasible without directly using knowledge on the document type. Furthermore, method implementation is specific for our database-internal representation of documents. For instance, it is a prerequisite that a flat object must not contain any other object.

```

(1) next(): logicalOID {
(2)   IF (SELF is a flat element) {
(3)     p := position of SELF within flat string;
(4)     p := SELF -> getPositionOfETAGO (p);
(5)     f := flat string of SELF;
(6)     p := f ->> find (p, STAGO)    //if the next element is contained in the same flat
                                     //database object, it begins after the end tag of the target element
                                     //‘find’ starts to search at byte offset identified by the first parameter.
                                     //It returns the byte offset where the second param. has been found, otherwise -1.
(7)     IF (p > -1) {    //next element is contained in the same flat
                           //database object, as its begin markup has been found
(8)       compute logical OID from p;
(9)       RETURN logical OID just computed; }; };
(10)    convert the (physical) OID of the next database object to logical OID;
        //trivial to identify next database object with structured representation
(11)    RETURN logical OID just computed; };

```

In the database, tag delimiters, e.g., STAGO, TAGC, are represented by special characters so that they cannot be mistaken with symbols ‘<’, ‘>’

within text, and search becomes more efficient.

Improving method performance. With OODBMSs, database method invocation is costly. Knowledge on the physical representation can be used to cut down the number of method invocations, and to reduce the parsing effort. In particular, it is worthwhile to avoid recursiveness. Consider the following implementation of method `getAll`.

```

(1) getAll (E: ElementType): {logicalOID} {
(2)   r := {};
(3)   IF (SELF is a flat element) {
(4)     p := position of SELF within flat string;
(5)     p_end := SELF -> getPositionOfETAGO (p);
(6)     f := flat string of SELF;
(7)     WHILE ((p < p_end) AND (p > -1)) {      //make sure that only elements within
                                                //target element are retrieved
(8)       p := f ->> find (p, concatenate (STAGO, TypeName (E)))
(9)       IF (p > -1)
(10)        IF ((isWithinBeginMarkup (p, f)) AND (E == type name of the element
                                                whose begin markup includes position p)) {
                                                //make sure that, e.g., AUTHORS is not found instead of AUTHOR
(11)          l := logical OID computed from p;
(12)          r := r ∪ {l}; }; }; }
(13)   ELSE
(14)     DO (children of SELF, element, e)    //iterate over the children of SELF
(15)       r := r ∪ (e -> getAll (E));
(16)   IF (ElementType (SELF) == E)
(17)     r := r ∪ {SELF};
(18)   RETURN r; };
```

On the contrary, a straightforward implementation would be recursive for all elements (as opposed to the one above that is only recursive for non-flat elements). Based on the MMF-DTD [S⁺94], we have conducted

experiments to verify that the first version is more efficient. If all elements are flat, and the root element is the target element of the original method invocation, the first version is faster by factor 1000 approximately. Naturally, the difference becomes smaller with fewer flat elements. If no elements are flat, the performance of the two versions is nearly identical.

4.2 Configurability Mechanisms

It is subject to configuration which elements are represented by individual database objects and which ones are flat. The configuration mechanisms are described next.

With our database application, documents of arbitrary type can be administered. Insertion of documents consists of the following steps:

1. The corresponding document-type definition is parsed. If the DTD is correct, a parser for instances of the DTD is generated. Furthermore, the DTD is - on a syntactical level - transformed to an SGML document that conforms to a specific DTD, the so-called *super-DTD*. The super-DTD is a document-type definition whose instances are document-type definitions. In the sequel, we will refer to any DTD different from the super-DTD as *application DTD*. For instance, the DTD from Figure 3 is an application DTD. It corresponds to the super-DTD-instance in Figure 7. (A fragment of) the super-DTD itself is contained in Figure 8.
2. At this point, the super-DTD instance contains exactly the information from the DTD. Attribute `elemName` of element type `ELEM`, to give an example, contains the element-type name, attribute `contentModel` contains the content model as a string. Furthermore, the super-DTD instances generated in Step 1 contain additional attributes that, initially, are instantiated with a default value. These attributes essentially contain information on the physical representation of element

```

<DOCTYPE docName=AGENDA ...>
<ELEM elemName=AGENDA ... contentModel='( HEADER , PROGRAMME )' ...>
    <ATTRIBUTE attrName=AUTHOR attrKeyDecl=CDATA attrKeyDef=IMPLIED ...>
</ELEM>

<ELEM elemName=HEADER ... contentModel='( LANGUAGE ?,( LOCATION |
ROOMNO ), DATE , INVITED )' ...></ELEM>

<ELEM elemName=PROGRAMME ... contentModel='( ITEM , ITEM *)' ...></ELEM>

```

Figure 7: Fragment of the Super-DTD instance corresponding to ‘Agenda’-
DTD

```

<!ELEMENT ELEM      (ATTRIBUTE*)>
<!ATTLIST ELEM      elemName      NAME  #REQUIRED
                    contentModel  CDATA  #IMPLIED ...>
<!ELEMENT ATTRIBUTE EMPTY>
<!ATTLIST ATTRIBUTE attrName      NAME  #REQUIRED ...> ...

```

Figure 8: Fragment of the Super-DTD

types or attribute types. For example, type **ELEM** has an attribute **FLAT**: value **NO** signifies that such elements are represented by individual database objects, **YES**, on the other hand, stands for a flat database-internal representation. By means of further attributes, the index structures are specified.⁴ - Summing up, in this step, i.e., Step 2, the physical representation of documents of a certain type is configured.

3. The document generated in Step 2 is parsed by a super-DTD-parser. In addition to checking the document's conformance to the DTD, the parser invokes database commands that generate the database objects that represent the document.
4. A database-internal bootstrap operation is invoked that, given the document inserted in Step 3, generates the corresponding database classes, index tables etc.
5. Now documents conforming to that application-DTD can be inserted into the database. The document parser that has been generated in Step 1 not only checks conformance to the DTD, but also invokes database operations generating the corresponding database representation, updates index tables etc.

Summing up, the physical representation of documents is configurable with element or attribute types being the granules of configurability. The dimensions of configurability are orthogonal to each other and transparent to the application programmer.

⁴Index structures can be turned on or off at a later stage by means of method invocations. The flat-/non-flat configuration, however, cannot be modified any more. Such a reorganization of the database would be extremely costly, and the need for such functionality has not yet arisen in our context.

4.3 The Transformation Algorithm from Documents' Logical Structure to their Physical Representation

In the sequel, we give the transformation algorithm that generates a document's database-internal representation from its logical structure. We will prove that the output of the transformation algorithm has certain important characteristics.

By definition, *Element type B is directly contained in element type A with regard to DTD D* if B occurs in the content model of A in D.

Definition 1 *Element type B is contained in element type A with regard to DTD D if*

1. *B is directly contained in A with regard to D, or*
2. *there is an element type C such that B is contained in C with regard to D, and C is contained in A with regard to D, or,*
3. *in D, A has an inclusion model that contains B, or*
4. *there is an element type C in D such that C has an inclusion model containing B, and A is contained in C with regard to D.*

The following lemma allows to derive information on the document type from a document that conforms to the underlying document-type definition. Due to the complex definition of containment on the type level, the lemma is not trivial. For instance, if an element a is directly contained in an element b, one cannot infer that the element type of a occurs in the content model of the element type of b (because of inclusions). From another perspective, the lemma shows that Definition 1 is meaningful.

Lemma 1 *If an element a of type A is directly contained in an element b of type B in some document of type D, then A is contained in B with regard to D.*

Proof. The proof is by induction on the depth of the document tree.

- a is directly contained in b , and b is the root of the document.

In this case, for a to be directly contained in b , either

1. A occurs in the content model of B , or
2. A occurs in the inclusion model of B .

In both cases, it follows directly from the definition that A is contained in B .

- a is directly contained in b . Furthermore, ‘ b contained in c . $\Rightarrow B$ contained in C .’.

For a to be directly contained in b ,

1. A occurs in the content model of B , or
2. A occurs in the inclusion model of B , or
3. A occurs in the inclusion model of an element type C , and there is an element c of type C such that a is indirectly contained in c .

It follows from Items 1, 3, 4 from the definition (corresponding to Items 1, 2, 3, respectively) that A is contained in B .

■

Lemma 2 *If an element a of type A is contained in an element b of type B in some document of type D , then A is contained in B with regard to D .*

Proof. The lemma immediately follows from Lemma 1 and Item 2 in the definition of ‘contains’.

■

In the transformation algorithm, the function `isFlat` with signature `isFlat (E: ElementType): BOOL` is used. It returns `TRUE` if E is contained in an element type that has been marked as flat in the corresponding super-DTD instance. In the algorithm, the document is traversed recursively in a

depth-first-like manner. If the type of the current element is not flat, a new database object is created and inserted into the tree structure that is already there (Lines 14-23). Otherwise, the current element's string representation is just appended to the current database object, which is flat (Lines 6-12).

The transformation algorithm is as follows:

```

(1) transform (e: Element, lastElementWasFlat: BOOL, currentObj: OLD, parentObj: OLD,
              root: BOOL): OLD {
(2)   IF (Type (e) -> isFlat()) {
(3)     IF (NOT (lastElementWasFlat)) {
(4)       currentObj := ELEMENT -> new();
(5)       insert currentObj as rightmost child of parentObj; };
(6)     IF (Type (e) is terminal element type) //e.g., CDATA, (#PCDATA)
(7)       currentObj -> append (textualContent (e))
(8)     ELSE {
(9)       { currentObj -> append (BeginMarkup (e));
(10)      DO (children of e, element, e') //iterate over the children of e, e' is loop var.
(11)        c := transform (e', TRUE, currentObj, parentObj, FALSE);
(12)      currentObj -> append (EndMarkup (e)); }; }
(13)   ELSE { //current element type is not flat
(14)     currentObj := ELEMENT -> new();
(15)     IF (NOT (root))
(16)       insert currentObj as rightmost child of parentObj;
(17)     store ElementTypeName (e) with currentObj;
(18)     store Attributes (e) with currentObj;
(19)     currentElementsFlat := FALSE;
(20)     c := NULL;
(21)     DO (children of e, element, e') //iterate over the children of e
(22)       { c := transform (e', currentElementsFlat, c, currentObj, FALSE);
(23)       currentElementsFlat := isFlat (Type (e')); }; }; }

```

ELEMENT is the database class described before, while Element is the type of SGML elements. - The initial invocation of transform is transform

(root, TRUE, NULL, NULL, TRUE). The actual implementation of the algorithm is non-recursive. Namely, an SGML parser has been extended to control the transformation that does not work recursively. Note that the database objects generated are untyped, i.e., they may either contain flat element types or represent non-flat elements. We say that *the database object is flat* or *the database object is non-flat*, respectively. By definition, an object becomes a flat one or a non-flat one by means of the assignments in Lines (7), (9), (12) or in Lines (17), (18), respectively. Thus, the definition of flat and non-flat database objects is an algorithmic one. From now on, this definition of flat database objects replaces the previous one.

The implementation of methods reflecting the SGML semantics such as `getAll` is based on the following lemmas.

Lemma 3 *After a type (i.e., either flat or non-flat) has been assigned to a database object, the type does not change any more in the course of the transformation algorithm.*

Proof. “ \Rightarrow ”: Consider a flat database object. The assignments making this object a non-flat one occur in Lines (17), (18). The object is generated immediately before (Line (14)). In consequence, it cannot happen that a flat object is subject to the assignments making it a non-flat one.

“ \Leftarrow ”: An object that is already non-flat cannot become a flat one later. Namely, non-flat objects are generated in Line (14) only. It can easily be seen that such an object does not become `currentObj` any more in the course of transformation after having specified that it is non-flat. ■

The following lemma shows that transformation by means of the algorithm is sound (cf. our remark on the implementation of `next`, and such knowledge has also been used for the implementation of `getAll` (Lines (3)-(12))).

Lemma 4 *A non-flat object is never contained in a flat one.*

Proof. Suppose a non-flat object was contained in a flat one. Then, there is a non-flat element e_1 that is directly contained in flat element e_2 . This requires that either

1. $\text{ElementType}(e_1)$ occurs in the content model of $\text{ElementType}(e_2)$,
or
2. there is an element type F s.t. $\text{ElementType}(e_1)$ occurs in the inclusion model of F , and there is an element e s.t. $\text{ElementType}(e) = F$ and e_1 is contained in e .

Case 1 cannot happen because $\text{ElementType}(e_1)$ would have to be a flat one. With regard to Case 2, it follows from Lemma 2 that $\text{ElementType}(e_2)$ is contained in F , and $\text{ElementType}(e_1)$ is contained in F . The last item from Definition 1 implies that $\text{ElementType}(e_1)$ is contained in $\text{ElementType}(e_2)$. This, however, is a contradiction to the definition of `isFlat`, because, in that case, $\text{ElementType}(e_1)$ would have to be flat. ■

A variant of the transformation algorithm is used in the context of document modification, i.e., in order to insert elements into documents that are already in the database.

5 DTD-Specific and Generic Declarative Access Mechanisms

By using SGML and database technology, we have come up with query mechanisms for a document collection characterized by the following features: (1) Formulating expressive queries is possible without knowing a query language, the DTD, and the underlying data model. (2) For different user groups, different mechanisms can be generated, closely matching the user group's needs. The description of these mechanisms and how to configure them covers a large part of this section. This query mechanism, though ex-

pressive, provides for a lower degree of expressiveness than others, as we will show. Our conclusion is to let the user choose between various, in our case three, query mechanisms differing with regard to expressiveness, but also with regard to intuitiveness and user-friendliness. In addition to query templates, there are an extension of the PAT algebra and OQL, together with methods from the schema. We will show that the extended PAT algebra is more expressive than query templates, and that OQL together with a relevant set of methods is more expressive than the extended PAT algebra. On the other hand, however, in order to formulate queries with the individual mechanisms, the user must have different levels of knowledge, as indicated in Table 1. More precisely, ‘n’ in the second column does not include the language of regular expressions and the underlying IRS, and ‘n’ in the third column does not reflect that the user has to understand the semantics of element- and attribute-type names.

	Modelling	(Syntax of the) query language	DTD
OQL	y	y	y
PAT	n	y	y
Templates	n	n	n

Table 1: Knowledge necessary to use different query mechanisms.

5.1 Query Templates

Query templates are automatically generated document-type-specific query forms. They may contain widgets of different types. It is subject to configuration which widgets are part of a query template. The following widgets are part of the framework.

- Entry field for element-content search. The figure is an example of such a widget, as seen in a WWW-browser. The user has to type in

a list of regular expressions, each of them separated by a blank space. The operation corresponding to the widget takes all elements of the respective type, in this case **SURNAME**. If AND is selected, it returns all documents containing elements of the type that contain all of the regular expressions. If OR is selected, it returns all documents containing elements of the type that contain one of the regular expressions.



- Entry field for IR search (information-retrieval search). An entry field for IR search actually consists of two fields, as can be seen in the figure. The user must type a concept to be searched for in the first entry field and a threshold value t in the second one. It must hold that $t \in [0;1)$. The corresponding operation takes all elements of the respective type. It returns all documents containing those elements that match the concept with a likelihood greater than the threshold value, as computed by the underlying IRS. - Instead of a concept, a query in the language of the underlying IRS can also be typed in.



- Entry field for attribute search. The corresponding operation takes all elements of the corresponding type, in this case **SECTION**. If AND is selected, it returns all documents containing elements whose value for attribute **SECQUAL** contains all of the regular expressions that have been typed in. If OR is selected, it returns all documents containing elements whose value for attribute **SECQUAL** contains one of the regular expressions that have been typed in.

- Entry field for structure search. The corresponding operation takes all elements of the first type, in this case **SURNAME**. For all such elements that are contained in one of the second type, in this case **AUTHOR**, and contains all of the regular expression that have been typed in, the corresponding document is returned, if AND has been selected. Analogously, with OR, only one regular expression must be contained.

- Entry field for search for physical characteristics. As opposed to the other atomic entry fields, these entry fields are hardcoded. However, they can be turned on or off by means of the configuration mechanisms. At this point, there is an entry field for document size allowing specification of a lower and upper bound.

The overall structure of a query template is depicted in Figure 9. The left column of widgets is for the document to be retrieved, the right column will be explained below. Results corresponding to individual entry fields in a column are combined using logical AND. Only those entry fields are considered where something has been entered.

More complex queries can be formulated using the widget for secondary structure search. Furthermore, one wants to specify documents by means of the (link) relationships that exist with other documents. In addition to those two columns of widgets, there is a pulldown menu with the following options: [NO LINKS], an element-type name/attribute-type name pair followed by a right arrow, and an element-type name/attribute-type name pair followed

by a left arrow. The semantics of the menu items is as follows:

- If [NO LINKS] is selected, the documents matching the entries in the left column are retrieved. Entities in the right column are ignored.



- If '-> <E>/<A>', e.g., '-> HYPLINK/REFERENC', is selected, selection is based on all pairs of documents (d_1 , d_2) such that d_1 matches the template entries in the left column, and d_2 matches the entries in the right column. The query returns all documents d_1 that contain an element of type E, this element has attribute A with value n, and n is the name of d_2 .



- If '<- <E>/<A>', e.g., '<- HYPLINK/REFERENC', is selected, selection is based on all pairs of documents (d_1 , d_2) such that d_1 matches the template entries in the left column, and d_2 matches the entries in the right column. The query returns all documents d_1 such that d_2 contains an element of type E, this element has attribute A with value n, and n is the name of d_1 .



The menu for secondary structure search is also subject to configuration.

The distinction between regular-expression search and IR search has been reflected by means of methods `hasTextualContentRegex` and `getIRSVValue` in Section 4. Analogously, query templates may contain both fields for element-content search and for IR search. As pointed out before, not only individual

concepts, but also complex query terms in the IRS-query language may be typed into entry fields for IR search. Consequently, it is not necessary to provide an AND-/OR-toggle for this widget type.

Specifying query templates. The DTD alone is not sufficient as a basis for automatic generation of query templates. Frequently, one wants to make available only a restricted set of types for declarative access. This corresponds to the notion of ‘view’ in the context of conventional database systems. The motivation why views should be part of the framework is as manifold as it is with view mechanisms in conventional systems. In principle, we see two alternative ways of specifying query templates.

1. The super-DTD is extended so that its instances contain the query-template specification. - Different ways of modelling the specification are conceivable. For example, there may be an additional element type `QUERYFORM` with attributes of type `IDREFS`. These references point to the different element- and attribute types to be included in the template. The type definition of `QUERYFORM` may be as follows:

```
...
<!ELEMENT QUERYFORM      EMPTY> ...
<!-- ATTLIST QUERYFORM
      CONTENTSEARCH      IDREFS
      ATTRSEARCH         IDREFS
      IR_SEARCH          IDREFS
-->
```

2. Each query template has a specification contained in a file.

With regard to Item 2, as one may need different templates for one document type, one may also want to freely add new views over time. However, it is important to ensure the consistency of the query-template specification

with the DTD. Otherwise, queries could be generated for which a solution cannot exist, and the user would not even notice it. But an operation which directly reads the specification from file and checks for its consistency would be too time-consuming with large DTDs. Hence, query-template generation must consist of two steps: First, the specification is read from file, its consistency to the DTD is checked, and it is inserted into the database. Then, a database-internal, consistent version of the specification can be accessed. - We have realized the first alternative and are now implementing the second one.

5.2 Other Declarative Access Mechanisms and a Comparison of their Expressive Power

An extension of the PAT algebra. The PAT algebra, originally described in [ST94], is a query language independent of the underlying data model. In our extension of the PAT algebra, query terms are generated by the grammar

```
e -> <Element-type name> |
    e UNION e |
    e INTERSECT e |
    e DIFF e |
    CONTENT_SELECT (e, r) |
    ATTR_SELECT (e, A, r) |
    IR_SEARCH (e, c, t) |
    e INCLUDES e |
    e INCL-IN e |
    REFERENCES (e, A, e) |
    REF-BY (e, A, e) |
    ID-REFER (e, A, e) |
    ID-REF-BY (e, A, e) |
    LB-SIZE (e, s) |
```

UB-SIZE (e, s) |
(e)

The term `<Element-type name>` stands for the set of all elements of the respective type. UNION, INTERSECT, and DIFF are set operators with the usual semantics. CONTENT_SELECT takes a set of elements and returns those where the content contains regular expression `r`. ATTR_SELECT takes a set of elements and returns those where attribute `A` contains regular expression `r`. IR_SEARCH takes a set of elements and returns those matching concept `c` (or the IR query `c`) with a probability greater than `t`, according to the underlying IRS. INCLUDES and INCL-IN take two sets of elements E_1 and E_2 and return the set of elements

$$E_1 \text{ INCL-IN } E_2 = \{e_1 \in E_1 \mid \exists e_2 \in E_2 \text{ s.t. } e_1 \text{ is contained in } e_2\}$$

$$E_1 \text{ INCLUDES } E_2 = \{e_1 \in E_1 \mid \exists e_2 \in E_2 \text{ s.t. } e_1 \text{ contains } e_2\}$$

REFERENCES, REF-BY, ID-REFER, and ID-REF-BY take two sets of elements E_1 and E_2 and return the set of elements

$$\text{REFERENCES}(E_1, A, E_2) = \{e_1 \in E_1 \mid \exists e_2 \in E_2 \text{ s.t. } e_1 \text{ has attribute } A \text{ with value } v, \text{ and } v \text{ is name of the document where } e_2 \text{ is contained in}\}$$

$$\text{REF-BY}(E_1, A, E_2) = \{e_1 \in E_1 \mid \exists e_2 \in E_2 \text{ s.t. } e_2 \text{ has attribute } A \text{ with value } v, \text{ and } v \text{ is name of the document where } e_1 \text{ is contained in}\}$$

$$\text{ID-REFER}(E_1, A, E_2) = \{e_1 \in E_1 \mid \exists e_2 \in E_2 \text{ s.t. } e_1, e_2 \text{ are contained in the same document, } e_2 \text{ has an attribute of type ID with value } v, e_1 \text{ has attribute } A \text{ of type IDREF}(S) \text{ containing } v\}$$

$$\text{ID-REF-BY}(E_1, A, E_2) = \{e_1 \in E_1 \mid \exists e_2 \in E_2 \text{ s.t. } e_1, e_2 \text{ are contained in the same document, } e_1 \text{ has an attribute of type ID with value } v, e_2 \text{ has attribute } A \text{ of type IDREF}(S) \text{ containing } v\}$$

LB-SIZE takes a set of elements and returns those whose size is greater than `s`, UB-SIZE returns those elements whose size is smaller than `s`.

The extensions, as compared to the original algebra [ST94], are the distinction between search on a syntactic level and IR search, the fact that

SURNAME	AND	OR		SURNAME	AND	OR	
<input type="text" value="Roth"/>	◇	◇		<input type="text" value=""/>	◇	◇	
LANGUAGE/ LANGQUAL	AND	OR		LANGUAGE/ LANGQUAL	AND	OR	
<input type="text" value=""/>	◇	◇		<input type="text" value="[Ee]nglish"/>	◇	◇	
SECTION / SECQUAL	AND	OR		SECTION / SECQUAL	AND	OR	
<input type="text" value=""/>	◇	◇		<input type="text" value="NEWS"/>	◇	◇	
ARTICLE	SW			ARTICLE	SW		
<input type="text" value=""/>	<input type="text" value=""/>			<input type="text" value=""/>	<input type="text" value=""/>		

<-HYPLINK/REFERENC

Figure 9: query template generated from the MMF-DTD

documents' secondary structure has been taken into account, and the fact that documents' physical characteristics have been considered.

OQL queries. The expressive power of OQL stems from the fact that methods from the database schema can be used within queries at liberty. The structure of an OQL-query is the same as with SQL. The select-clauses specifies what is to be selected. The from-clause specifies which database classes, or, more generally, which sets, the query refers to. The where-clause contains a condition that must be fulfilled by the query result. All variables occurring in the query must be bound in the from clause. - The reader is referred to [Cat94] for more information on OQL.

Illustrations. For illustration purposes, consider the query template in Figure 9. The template entries correspond to the query "Select all documents containing an element of type **SURNAME** whose textual content contains 'Roth', and that are referenced by a document containing an element of type **SECTION** whose value of attribute **SECQUAL** contains 'NEWS' and containing an element of type **LANGUAGE** whose value of attribute **LANGQUAL** contains 'English' or 'english'.". The corresponding extended PAT expression is

REF-BY (MMF INCLUDES CONTENT_SELECT (SURNAME, 'Roth'),


```

REFERENC,
HYPLINK INCL-IN (MMF INCLUDES ATTR_SELECT (SECTION, SECQUAL, 'NEWS'))
INCLUDES ATTR_SELECT (LANGUAGE, LANGQUAL, '[Ee]nglish'))

```

The corresponding OQL expression is

```

select D0.name
from D0 in Document, D1 in Document,
    P0 in D0.root -> getAll ('SURNAME'),
    P1 in D1.root -> getAll ('SECTION'),
    P2 in D1.root -> getAll ('LANGUAGE'),
    P3 in D1.root -> getAll ('HYPLINK')
where P0 -> hasTextualContentRegex ('Roth') and
    P1 -> hasAttrValue ('SECQUAL', 'NEWS') and
    P2 -> hasAttrValue ('LANGQUAL', '[Ee]nglish') and
    P3 -> getAttrValue ('REFERENC') = D0.name

```

The following lemmas reflect the expressive power of the different mechanisms.

Lemma 5 *The extended PAT language is more expressive than query-template entries.*

Proof. The proof is by defining a mapping from query-template entries to expressions in the extended PAT language. The full mapping is given in [B97]. To illustrate the mapping, consider the widget for attribute search. Let **E** be the respective element-type name, and **A** be the attribute-type name. With $r_1 \dots r_n$ being the input to the respective field, the corresponding expression is `<root-element-type> INCLUDES (ATTR_SELECT (E, A, r_1) $\pi \dots \pi$ ATTR_SELECT (E, A, r_n))` with $\pi \in \{\text{UNION, INTERSECT}\}$. In the opposite direction, it is obvious that, e.g., the extended PAT expression `A INCLUDES B INCLUDES C` cannot be mapped to any query-template input. ■

Lemma 6 *The extended PAT-query language is less expressive than OQL, together with the methods given in Section 4.*

Proof. The proof is by defining a mapping of extended PAT expressions to OQL statements. The proof is recursive over the structure of query-algebra terms. Again the full mapping is contained in [B97]. As an example, let Q be the OQL query corresponding to the PAT expression e . Then `ATTR_SELECT` (e , A , r) is mapped to

```
select p from p in Q where (p -> hasAttrValueRegex ('<A>', '<r>'))
```

In the opposite direction, it is obvious that, e.g., the OQL query

```
select p, p -> getFirst ('CHRNAME')
from e in ElementType, p in e -> getElements()
where (e.name == 'AUTHOR')
```

cannot be mapped to any extended PAT expression. ■

For evaluation, both input to query templates and extended PAT expressions are mapped to OQL expressions. - Declarative access mechanisms are also relevant in the following section.

6 Hyperlink Rendition Mechanisms in HyperStorM

With WWW-based access to the database application, documents can be converted to HTML. Conversion is specified by means of a style-sheet contained in a file.

Hyperlink rendition. Documents may contain references to other documents. Usually, such references are made explicit within the document with hypertext anchors. With our system, rendition of anchors of links pointing to other documents in the database may depend on characteristics of the documents referenced. As a special case of such rendition, only anchors of sound links are converted to HTML anchors to avoid some cases of navigation failure. The layout specification specifies how link anchors are encoded in documents of the respective type.

In Section 4, it has been described how the physical representation of documents and document components can be configured using the super-DTD. These mechanisms, however, are not used to specify document conversion for the following reasons:

- An initial configuration must have been specified before documents are inserted into the database. This does not have to be the case for the conversion specification. Furthermore, a higher degree of flexibility and ease of modification is necessary with the conversion specification, as compared to the configuration specification. It is appropriate if the configuration is altered by means of method invocations, but this is too complicated and inflexible for the conversion specification.
- From an organizational perspective, while the database-internal configuration should be specified by the DTD designer, this is not necessary for document rendition, as readers' individual preferences may be reflected.
- The super-DTD has been designed to represent information on individual (element- or attribute-) types. But the super-DTD instance would become too big if information from several style-sheets was included. There should be no restrictions to the number of style-sheets.

Incorporating information on documents referenced into hyperlink anchors. With our system, information on documents directly or indirectly referenced can be used to render the corresponding hyperlink anchor in the referencing document in a very flexible way. The core idea is that, for an element type whose instances are hyperlink anchors, the style-sheet contains a database query. The query specifies the information to be included in the hyperlink anchor. For this purpose, the full expressive power of OQL can be exploited.

In this context, there are two problems impeding why OQL queries cannot just be written down and executed during document conversion.

1. Query results must be of a type that can be displayed within a HTML document.
2. Within such queries, one would like to refer to the particular hyperlink anchor. For example, one would like to formulate the query “Select all elements of type `AUTHOR` within the document that is referenced by the hyperlink whose anchor is currently being rendered.”. So far, there is no straightforward way to refer to the hyperlink whose anchor is currently being rendered in a query.

The solution to Item 1 is as follows: the query is parsed. Then, the query-result type is looked up. If the type is not available for display, e.g., because the query result is a set of database OIDs or a set of instances of bulk types, the query in the style-sheet is ignored. The current version of the system accepts only sets of strings.

With regard to Item 2, OQL-queries within the style-sheet may contain the symbols `ANCHORELEM` and `LINKATTR`. Before query evaluation, the symbol `ANCHORELEM` is replaced by an expression that can be interpreted by the query processor. In a nutshell, it is the logical OID of the anchor element. Since the replacement is carried out within the database, it is ensured that the new query is correct. Consider the following example:

```
select  p -> getElementText()
from    d in Document, p in d.root -> getAll ('LANGUAGE')
where   (d.name = (ANCHORELEM -> getAttrValue (LINKATTR)))
```

The query selects all `LANGUAGE`-elements in the document referenced (cf. Figure 1).

Expanding instances of element types is meaningful for informational and non-structural element types. For instance, one may want to display all

instances of a non-structural element type `KEYWORD` in the referencing document. The usefulness of such an expansion for structural element types is limited, except when these element types are at the same time informational, such as, say, `TITLE`.

Summing up, the drawbacks pointed out in the introduction are alleviated as follows: checking for a link target before activating the anchor prevents navigation failure. Displaying meta-information alleviates the “Expectations vs. actual Content”-problem and is a measure against “getting lost in hyperspace”.

With the functionality described so far, document conversion should be part of the services offered by the database. It would be an unnecessary step to generate SGML from the database content first and then transform the result to HTML. Furthermore, conversion on the client would be inefficient because of unnecessary communication over the network to obtain information on the documents referenced. Besides that, and most importantly, materialized views are used to avoid executing queries during document conversion. Their consistency is ensured by the database system.

Using materialized views. It is advantageous to materialize the conversion result within the database. In principle, an arbitrary number of different views on the same document being the result of different conversion parameters is conceivable. To cope with the requirement of materializing more than one view, but being restricted to a limited number of materializations, we have implemented a simple page replacement strategy (LRU). With our system, documents within the database may be modified. Then, materialized views have to be updated after the corresponding documents have been altered. The problem is aggravated by the fact that updates of individual documents cannot be seen in isolation. Instead, documents that directly or indirectly reference the modified document also have to be taken

into account. In this context, we assume that each document can reference any other document. Assertions that would allow restrictions of this assumption would be helpful, but cannot be made in the general case. With our system, a view is updated when a document is accessed, if any document has been modified after the view's last generation date. With frequent updates, a conceivable refinement is to differentiate between modifications having an impact on conversion of referencing documents and modifications without such effects. In this case, however, conversion specifications have to be administered by the database.

7 Conclusions

One objective of this article was to point out how the combination of database technology and SGML can be exploited in order to ease access to the document collection, both declaratively as well as by means of navigation. In order to ease declarative access, a mechanism to express information needs has been designed that is expressive, while, at the same time, neither knowledge of a query language nor the document type or the underlying modelling are needed to use it. This so-called query-template mechanism can be configured to match the needs of different user groups. We exploit the fact that SGML documents conform to a document-type definition. Such a document-type definition is comparable to a database schema in that both can be seen as integrity constraints on the data. HTML would not be very useful in this particular context due to its high genericity. The HTML-DTD does not impose any real constraints on documents' logical structure. - The query-template mechanism, though expressive, is less expressive than other query languages. Thus, in order to provide the user with a choice of declarative access mechanisms, queries can also be formulated using two other languages. These languages have also been described, and their expressiveness has been compared.

Navigational access is eased as follows: meta-information on the documents can be exploited for hyperlink rendition. Database technology is advantageous to speed up conversion, and to ensure consistency of conversion outputs that have been materialized. The approach would not work well with HTML because, with previous versions of HTML, no mechanisms are provided to model meta-information, and with HTML 3.2 it cannot be ensured that meta-information is modeled consistently. It is left to the author of the document which meta-information is provided, and how he provides it.

With the description of documents' database-internal representation, we have covered several original features. With our hybrid representation of documents, elements can or cannot be represented by individual database objects. The actual representation is subject to configuration. The initial configuration is specified by means of an SGML document. To recognize the logical structure of physically unstructured document components, it is not necessary to use the DTD at this stage. We have come up with an object-oriented modeling of structured documents, together with methods that are the basis for declarative access. Elements have a (logical) OID even though they do not have to be represented by an individual database object. With method implementation, it has been reflected that, with OODBMSs, method invocations are expensive.

In the future, we wish to evaluate the behavior of a real user community. We will examine the expressive power of query templates and see whether new primitives are needed. - In previous work, we have examined how to reflect the semantics of HyTime architectural forms in the database and have come up with an implementation [BA94, BAK95]. The HyTime-link model is more sophisticated than the one that has been considered so far. Arbitrary document portions can be referenced, independent of the structure that has been made explicit with SGML markup. It may be worthwhile to extend

the reflections with regard to link conversion to the HyTime model.

References

- [ACM93] S. Abiteboul, S. Cluet, and T. Milo. Querying and Updating the File. In R. Agrawal, S. Baker, and D. Bell, editors, *Proceedings of the International Conference on Very Large Data Bases*, pages 73–84. VLDB Endowment, 1993. Dublin, Ireland.
- [ACM95] S. Abiteboul, S. Cluet, and T. Milo. A Database Interface for File Update. In *Proceedings ACM SIGMOD*, pages 386–397. ACM Press, 1995.
- [AKM95] K. Andrews, F. Kappe, and H. Maurer. The Hyper-G Network Information System. *J.UCS*, 1(4):206–220, 1995.
- [B⁺94] G.E. Blake et al. Text / Relational Database Management Systems: Harmonizing SQL and SGML. In *Proceedings of the First International Conference on Applications of Databases*. Lecture Notes in Computer Science, Springer Verlag, June 1994.
- [Bö97] Klemens Böhm. *Using Object-oriented Database Technology for Structured Document Storage (in German, forthcoming)*. PhD thesis, Technical University of Darmstadt, 1997.
- [BA94] Klemens Böhm and Karl Aberer. Storing HyTime Documents in an Object-Oriented Database. In Nabil R. Adam, Bharat Bhargava, and Yelena Yesha, editors, *Proceedings of the Third International Conference on Information and Knowledge Management*, pages 26–33. ACM Press, November 1994.

- [BAK95] Klemens Böhm, Karl Aberer, and Wolfgang Klas. Building a Configurable Database Application for Structured Documents, 1995. Accepted for publication in *Multimedia - Tools and Applications*.
- [C⁺94] V. Christophides et al. From Structured Documents to Novel Query Facilities. In *Proceedings ACM SIGMOD*. ACM Press, May 1994.
- [Cat94] R.G.G. Cattell, editor. *The Object Database Standard: ODMG-93*. Morgan Kaufmann Publishers, 1994.
- [CCH92] J.P. Callan, W.B. Croft, and S.M. Hardig. The INQUERY Retrieval System. In *Proceedings of the Third International Conference on Database and Expert Systems Application*, pages 78–83. Springer Verlag, 1992.
- [Con87] J. Conklin. Hypertext: An Introduction and Survey. *IEEE Computer Magazine*, pages 17–41, September 1987.
- [Fre] E.D. Freese. The Transformation of SGML Documents for Presentation on the World Wide Web. <http://www.sil.org/sgml/freese.html>.
- [HHM94] A. Haake, C. Hüser, and W. Möhr. Milestone M1; BERKOM Project:CLIP-ING; Workpackage 2: System Architecture. Technical report, GMD-IPSI, November 1994.
- [HTM] HyperText Markup Language (HTML). Available under "http://www.w3.org/pub/WWW/MarkUp/".
- [ISO86] Information Technology - Text and Office Systems - Standardized Generalized Markup Language (SGML), 1986.

- [ISO96] Document Style Semantics and Specification Language (DSSSL), 1996.
- [kbo94] Guidelines for Electronic Text Encoding and Interchange, April 1994.
- [KSS95] V. Kashyap, K. Shah, and A. Sheth. *Multimedia Database Systems: Issues and Research Directions*, chapter Metadata for building the MultiMedia Patch Quilt. Springer Verlag, 1995.
- [Mac91] I.A. Macleod. A Query Language for Retrieving Information from Hierarchic Text Structures. *The Computer Journal*, 34(3):254–264, 1991.
- [MMM96] A. Mendelzon, G. Mihaila, and T. Milo. Querying the World Wide Web, 1996. PDIS’96.
- [NBY95] G. Navarro and R. Baeza-Yates. A language for queries on structure and contents of textual databases. In *Proceedings of 18th ACM Conference on Research and Development in Information Retrieval (SIGIR’95)*, July 1995. Seattle, WA, U.S.A.
- [O⁺95] M.T. Özsu et al. An Object-Oriented Multimedia Database System for a News-on-Demand Application. *Multimedia Systems*, (3):182–203, 1995.
- [QRS⁺95] Dallan Quass, Anand Rajaraman, Yehoshua Sagiv, Jeffrey D. Ullman, and Jennifer Widom. Querying Semistructured Heterogeneous Information. In *Proceedings of the Fourth International Conference on Deductive and Object-Oriented Databases*, pages 319–344, 1995. Singapore.
- [RMW94] Martin Röscheisen, Christian Mogensen, and Terry Winograd. Shared Web Annotations as a Platform for Third-Party Value-

Added Information Providers: Architecture, Protocols, and Usage Examples. Technical Report STAN-CS-TR-97-1582, Stanford University, November 1994.

- [S⁺94] Klaus Süllow et al. MultiMedia Forum - an Interactive Online Journal. In Christoph Hüser, Wiebke Möhr, and Vincent Quint, editors, *Proceedings of Conference on Electronic Publishing*, pages 413–422. John Wiley & Sons, Ltd., April 1994.
- [SDAMK95] R. Sacks-Davis, T. Arnold-Moore, and A. Kent. A Standards Based Approach to Combining Information Retrieval and Database Functionality. *International Journal of Information Technology, World Scientific*, 1(1):1–16, 1995.
- [SM83] Gerard Salton and Michael J. McGill. *Introduction to Modern Information Retrieval*. McGraw-Hill Book Company, first edition, 1983.
- [ST94] A. Salminen and F.W. Tompa. PAT Expressions: an Algebra for Text Search. *Acta Linguistica Hungarica*, 41(1):277–306, 1994.
- [Sto96] Michael Stonebraker. *Object-Relational DBMSs*. Morgan Kaufmann Publishers, Inc., 1996.
- [VAB96] M. Volz, K. Aberer, and K. Böhm. Applying a Flexible OODBMS-IRS-Coupling to Structured Document Handling. In *Proceedings of the 12th International Conference on Data Engineering*, pages 10–19, 1996. New Orleans.
- [VML95] VODAK V 4.0 User Manual. Technical Report 910, GMD-IPSI, April 1995. St. Augustin.

- [YA94] T.W. Yan and J. Annevelink. Integrating a Structured Text Retrieval System with an Object-Oriented Database System. In *Proceedings of the International Conference on Very Large Data Bases*. VLDB Endowment, 1994. Santiago, Chile.