

# TUM

INSTITUT FÜR INFORMATIK

## Indexing XML as a Multidimensional Problem

Michael G. Bauer, Frank Ramsak, Rudolf Bayer



TUM-I0203

Mai 02

TECHNISCHE UNIVERSITÄT MÜNCHEN

TUM-INFO-05-I0203-0/1.-FI

Alle Rechte vorbehalten

Nachdruck auch auszugsweise verboten

©2002

Druck:            Institut für Informatik der  
                  Technischen Universität München

# Indexing XML as a Multidimensional Problem

Michael G. Bauer

Frank Ramsak

Rudolf Bayer

31st May 2002

## 1 Introduction

XML is widely seen as the lingua franca of the Internet. Originally designed to become the successor of HTML, XML has found its way into many unexpected parts of applications, ranging from simple formats for data exchange to archiving data in XML. With the growing need to deal with large collections of XML documents as well as with the rapid increase in the document sizes there is a strong demand to store and query XML data in databases. This ranges from the development of new, efficient index structures for XML to mapping schemes for XML to relational and object-oriented database systems. As relational database management systems (RDBMS) hold the largest market share there was and still is intensive research going on to efficiently store XML in these systems. Several mapping schemes have been proposed [5, 4] in the literature but to our knowledge there has never been an extensive analysis of a mapping scheme in conjunction with index structures.

In our work we discuss a multidimensional approach for indexing XML. We propose a multidimensional mapping scheme for XML to relational DBMS and discuss the performance of UB-Trees and compound B-Trees for the indexing of this mapping scheme.

The rest of the paper is structured as follows. At first we motivate the problem of XML indexing and present a modelling of the XML document in a multidimensional universe (Section 2). We explain the MHC (multidimensional hierarchical clustering) technique (Section 3) and continue by presenting an implementation of our multidimensional approach by using MHC, the UB-Tree and various compound B-Trees as multidimensional indexstructures (Section 4). We also present some performance experiments in Section 5. We conclude the paper with a short summary in Section 6.

## 2 The Problem: Storing XML in RDBMSs

Many approaches have been made to store XML in relational database systems. All approaches use a similar concept though. First the XML document is split into parts of a previously defined granularity. These parts are stored in the RDBMS. Queries on the XML documents which are written in an XML query language have to be rewritten to SQL before being processed by the RDBMS. The results of the SQL queries are documents. The projection is either done via methods like XSLT or the projection results are assembled directly from the database. Our approach especially takes care of both cases.

Besides storing XML data, querying large amounts of XML data is a challenging problem. Due to its graph-like nature the classical set oriented query languages in general are not powerful enough. After several proposals for query languages (mainly from the fast evolving field of semistructured data) the W3C started a working group to formulate a query language for XML. The current proposal XQuery though is not yet a recommendation of the W3C.

We have identified two fundamental parts of a query for XML. Consider the following query in Pseudo-SQL which should retrieve the value of a tag in a document containing the tag ABC: **select <tag1> from xmlbase where tag2='ABC'**; This query can be (similarly to queries in the relational case) split up into two steps. The selection part identifies the document(s) which contain(s) the pattern matching a predicate  $s$ . The projection part in contrast returns only those

part(s) of the document(s) which are stated right after the **select** statement as a set of paths. The mentioned problems can be defined more formally:

**Definition 1 (selection problem)** *Let  $X$  be a set of XML documents stored in a database where each document is described by a persistent unique identifier  $Id$ . The selection problem is defined as returning  $Id$  for those documents where the predicate  $s$  from the query is evaluated to true.*

**Definition 2 (projection problem)** *Let  $I$  be a set of persistent unique identifiers and  $L$  a list of paths in a projection list. The projection problem is defined as returning the content of those paths from  $L$  from the documents referenced by  $I$ .*

When talking about the projection problem we sometimes use the term "reconstruction of (parts of) the document". We refer to the fact that the desired document is completely or partially assembled from the database into its original state.

In our approach we do not tackle only one of the above mentioned problems (either selection or projection) but both.

One has to be aware of the fact that the projection in the case of XML documents is fundamentally different from the relational world. The operations in the relational algebra deal with tuples and sets as the basic units. Tuples in RDBMSs are typically small compared to the size of an XML document, therefore tuples are normally retrieved as a whole during the selection process and the projection always processes the tuples that resulted from the selection. For efficiently answering queries in the relational world it is therefore sufficient to only speed-up the selection. For XML documents the scenario is slightly different as the granularity of an XML document can be seen on different levels. A very rough granularity is based on documents (which are identified by a persistent unique identifier). Tags as the building block of XML documents are another level of granularity, a very fine granularity would be based on words or letters. Depending on the storage of XML in the RDBMS the projection works on a completely different position of the XML document than the selection. The consequence is that an index that is suitable for the selection is rarely suitable to speed up projection as well. We will see later that selection and projection can even be orthogonal and require very specialized index support.

## 2.1 XML - A Multidimensional Model

Paths are a fundamental feature of XML. Many approaches to speed up queries therefore concentrated on the efficient indexing of paths and path expressions. When discussing paths it is also important to note the order of paths which is inherent in the XML documents. The ordering of XML documents is a very important aspect especially in the case of indexing. The DTD of an XML document already defines the ordering of the tags in the document. This is especially important when tags with the same substructure occur several times but with different data.

### Example 1

```
<author><FN>Michael</FN><LN>Bauer</LN></author>
<author><FN>Rudolf</FN><LN>Bayer</LN></author>
```

Order is also important for query languages. When the above example (a fragment from a larger document) is queried with the predicate **author/FN = Michael and author/LN = Bayer** the semantics of the predicate is not clear at first hand. In the document fragment the paths to the values are the same but the structure and the order of the paths is of great importance as it acts as a grouping feature that is very relevant for queries. The above predicate can be reformulated so that the semantics is clear and the two parts of the query have to be valid in the same subtree. A correctly rewritten version of the above predicate in XPath so that the document which contains the above fragment is returned is formulated as **author[FN = "Michael" and LN = "Bayer"]**. To be able to evaluate such a predicate it is apparently necessary to preserve the document structure for the stored document.

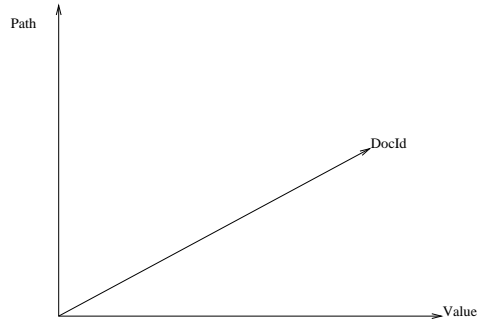


Figure 1: XML documents in a three dimensional cube

Finally ordering might have a severe impact on the performance of answering queries. In every database instance the data is stored in a certain physical order on external storage devices. Clustering the stored data in document order usually can be achieved, nevertheless an order independent of the document order might be more useful for certain queries. Nevertheless, the original order of the stored document has to be preserved in some way, otherwise it is impossible to reconstruct the document in the same order as it was originally available.

In the following we propose three building blocks of XML documents.

**Paths:** The notation of paths is a basic concept of XML and query languages for XML. Due to the graph-like nature of XML it is necessary to preserve path information and order for query processing.

**Values:** We define values as the content of XML tags. Our proposed scheme can deal with both, data-centric and document-centric XML. Attributes in XML are modelled by using the @notation in the paths as known from XPath.

**Document Identifiers:** Document identifiers group paths for one document and are the results in the above mentioned selection. We assume that this identifier is available from the XML data itself. If this is not the case it can be easily computed when the document is processed before it is inserted into the database.

Summing up our results we can define a set of XML documents which are stored in our database as follows.

**Definition 3** Let  $P$  be a set of paths,  $\langle_{path}$  the order of the paths,  $V$  a set of values and  $id$  a document identifier.

$$XMLdocs = \bigcup \{(P, \langle_{path}, V, id) | id = docid\}$$

For better visualization the three dimensions can be presented as a cube (Figure 1). In this three dimensional model we can now easily answer queries for both the selection and the projection problem as follows.

Both selection and projection restrict the threedimensional universe in two dimensions. The input for evaluation of the XPath predicate of the selection is one (or more) path expressions and one (or more) values which correspond to the path expression and form a predicate. The output of the selection is a set of document identifiers which match path expression(s) and the value(s).

In the projection the document identifiers and the path expressions are the input of the query. The results are values. For better readability and post-processing capabilities the document identifiers as well as the output paths are sometimes returned as well.

### 3 MHC

MHC (Multidimensional Hierarchical Clustering) is a technique that was originally developed for data warehousing applications [7]. In data warehousing MHC is used to cluster data with respect to multiple hierarchical dimensions. We briefly describe the technique in an overview, show its application in our context of XML indexing, how we use MHC to preserve order and how to store paths in a compact form.

It is well known by now that XML documents form a hierarchy with nodes and edges. We assume that each XML document has a maximal hierarchy of depth  $h$  leading to an overall of  $h + 1$  levels in each document. We use a slightly different tree representation as it is usually common. Identical paths which occur several times within a subtree are rewritten by using repetition numbers to preserve uniqueness and order. The paths from example 1 are rewritten to  $\langle \mathbf{author} \rangle[1]\langle \mathbf{FN} \rangle[1]$ ,  $\langle \mathbf{author} \rangle[1]\langle \mathbf{LN} \rangle[1]$ ,  $\langle \mathbf{author} \rangle[2]\langle \mathbf{FN} \rangle[1]$ ,  $\langle \mathbf{author} \rangle[2]\langle \mathbf{LN} \rangle[1]$ .

In the following we treat the repetition numbers as an own hierarchy level. The set of levels is ordered according to the document structure. Each hierarchy level  $i$  is a set of sets  $L$  ( $L = \bigcup_{i=0}^n L_i$ ), where each set  $L_i$  consists of nodes  $m_k^i$ .  $L_0$  is defined to be the root level. Due to the hierarchical nature every member  $m_k^i$  has a (varying) number of children. A function  $ord_m$  defines a numbering scheme for the children of  $m_k^i$  and assigns each child of  $m_k^i$  a number between 1 and the total number of children of  $m_k^i$ , i.e.  $ord_m : children(m) \rightarrow \{1, \dots, |children(m)|\}$ . The function  $ord_m$  needs to be order preserving so that nodes are sequentially numbered according to document order. We call each element of  $\{1, \dots, |children(m)|\}$  a surrogate of a node. To encode paths through the hierarchy we introduce the concept of compound surrogates. A compound surrogate is formed by recursively concatenating the compound surrogate of the father node with the surrogate of the current node. More formally the compound surrogate for a node is defined as follows:

$$cs(m^i) = \begin{cases} ord_{father(m^i)}(m^i), & \text{if } i = 1 \\ cs(father(m^i)) \circ ord_{father(m^i)}(m^i), & \text{otherwise} \end{cases}$$

**Example 2** Using figure 2 the path  $\langle \mathbf{author} \rangle[1]\langle \mathbf{FN} \rangle[1]$  is transformed into the compound surrogate 2.1.1.1

Compound surrogates can be very efficiently stored in a compact binary representation. The upper limit for each surrogate can be calculated by the formula

$$surrogate(i) = max\{cardinality(children(m)) \text{ where } m \in level(i - 1)\}$$

The length of the surrogates can be chosen sufficiently large in advance and although the fixed length leads to static boundaries of the surrogates this can be neglected. Extending every surrogate by one bit doubles the number of children that can be addressed at every level.

MHC enables us to find a compact and order preserving representation for paths in XML documents. It also fixes the problem of long equal prefixes which has already been addressed in work on special index structures for XML [4]. It is also possible to evaluate path expression on compound surrogates. Evaluating simple path expressions (i.e. full qualified paths) in MHC is equivalent to point queries on the compound surrogates while other path expressions are equivalent to range queries or combinations of point and range queries. In a separate forthcoming paper we asserted that all 13 location steps which are defined in XPath 1.0 and XPath 2.0 can be implemented using the representation of paths as compound surrogates.

### 4 Implementation

Over the last year several new indexes for XML documents were developed. There were also several proposals for mapping schemes of XML documents to relational database management systems. To our knowledge there was no thorough comparison of a mapping scheme with different

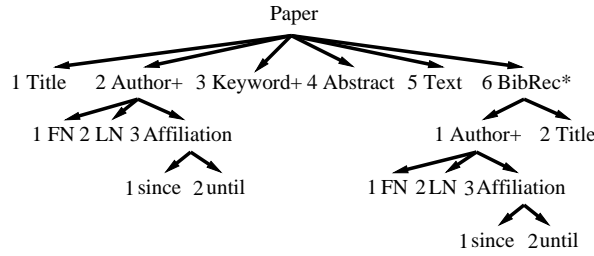


Figure 2: XML document hierarchy and MHC surrogates

did	val	surr
1	Rudolf	0010 0001 0001 0001 0000 0000 0000 0000
1	Bayer	0010 0001 0010 0001 0000 0000 0000 0000

Table 1: Relation `xmltriple`

indexing methods instead indexes were often chosen adhoc and without further discussion of the effects.

The implementation of our mapping scheme is based on two relations. The core is a table with three attributes, which we refer to as `xmltriple`. `xmltriple` holds the attributes `did`, `surr`, and `value` (see Table 1 for an example with two tuples and 4 bits per surrogate).

For mapping XML document paths to compound surrogates we use an additional table `typedim` with the two attributes `path` and `surr` (Table 2).

For our measurements we compared three different indexing methods for our proposed mapping scheme. Since we claim that XML indexing is a multidimensional problem we chose the UB-Tree (threedimensional) and two variants of B-Tree compound indexes for the `xmltriple` relation.

## 4.1 The Data

For our measurements we have chosen a DTD for scientific papers. This DTD represents a typical application in the field of digital libraries. We generated 10000 different documents using the XMLGenerator tool [8]. The raw size of the XML data is approx. 50 MB. We have used a similar distribution [4] as in the DBLP database [9] for authors which results in approx. 400 different authors for 10000 documents. From these documents we generated flat files for bulk loading the databases. The sizes of the database are approx. 25 MB and they differ slightly depending on the used index.

## 4.2 The Query

Our data is a typical library application. In accordance with this application we also chose a typical query.

### Example 3

```
select titel, keywords
```

surr	path
0010 0001 0001 0001 0000 0000 0000 0000	/Author[1]/FN[1]
0010 0001 0010 0001 0000 0000 0000 0000	/Author[1]/LN[1]

Table 2: Relation `typedim`

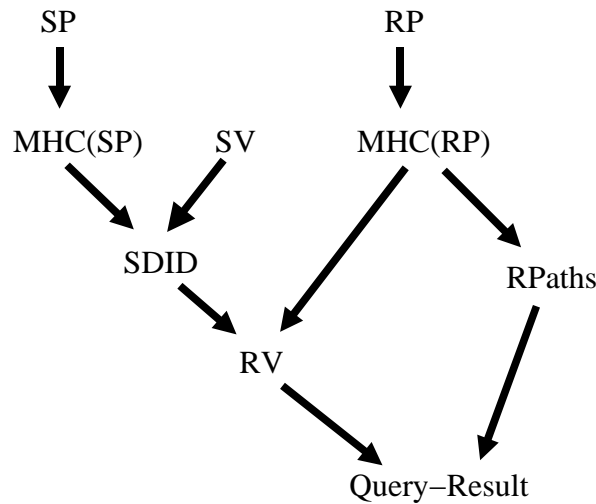


Figure 3: Steps to rewrite the query into SQL

```

from xmldata
where /author[FN = "Michael" and LN = "Bauer"]

```

The query returns the title and keywords of all scientific papers that were written by a certain author.

To be processed by the RDBMS the above query has to be rewritten to SQL by using the method of figure 3.

The initial query is separated into the already mentioned two components *selection* and *projection*. The selection consists of a selection path (SP) and selection value (SV), which can be seen in the left branches of the diagram in figure 3. The paths are mapped to compound surrogates by a lookup in the **typedim** table (MHC(SP)). The compound surrogates and search values are then used to identify the document ids which satisfy the selection predicate (SDID) by querying the **xmltriple** relation.

The projection, which outputs the result values (RV), is processed similarly. The result paths (RP) are transformed into compound surrogates. In the next step the result values are returned by using the retrieved document ids and the compound surrogates as input. More complex queries might require recursive application of this schema.

As we are examining a relational mapping of XML data the queries of each step have to be rewritten to SQL statements. Some of the statements depend on the output of previous queries. This may lead to long statements (especially for many hits in the document ids) and for ease of presentation we only give very short example statements.

**Example 4** 1. Step: *MHC(SP)*

```

select surr from typedim
  where attr like 'Author[%]/FN[1]' order by surr;
select surr from typedim
  where attr like 'Author[%]/LN[1]' order by surr;

```

2. Step: *SDID*

```

select distinct x1.did from xmltriple x1, xmltriple x2
  where
    (x1.surr = 0b00100001000100010000000000000000

```



```

or x1.surr = 0b00100010000100010000000000000000
.....      result of MHC(SP)
or x1.surr = 0b00101000000100010000000000000000
and
x1.val = 245650
and
(x2.surr = 0b00100001001000010000000000000000
or x2.surr = 0b00100010001000010000000000000000
or .....      result of MHC(SP)
or x2.surr = 0b00101000001000010000000000000000)
and
x2.val = 813220
and
x1.did = x2.did
order by x1.did

```

### 3. Step: MHC(RP)

```

select surr from typedim
  where attr like 'Title[1]'
        or attr like 'Keyword[%]' order by surr;

```

### 4. Step: RV

```

select did,val,attr,typedim.surr from xmltriple,typedim
  where (typedim.surr = 0b00010001000000000000000000000000
or typedim.surr = 0b00110001000000000000000000000000
or .....      result of MHC(RP)
or typedim.surr = 0b00111000000000000000000000000000)
and ( did = 76 or .....      result of SDid
or did = 9783 )
and typedim.surr = xmltriple.surr
order by did,typedim.surr

```

## 5 Measurements

For our measurements we ran the above mentioned query on the different indexed tables. All measurements were performed with the relational database system TransBase<sup>1</sup>, which won the European Information Technology Prize 2001 for its pioneering implementation of UB-Tree indexes. We chose a page size of 2KB and limited the database cache to 128 KB. With a cache size this small not many pages can be kept in the cache. Completely eliminating the cache would severely decrease performance as even index pages would no longer reside in the cache.

The database system was installed on a Sun Ultra 10 (400MHz, 512MB main memory) and the benchmarks were performed on a Seagate ST39111A (73.4GB) U160 hard disk.

In the following we use abbreviations to denote the different indexing methods. *xmlub* denotes the indexing with the threedimensional UB-Tree, the index attributes are *did*, *surr*, *value*. *xml-didsurr* denotes indexing with a compound B-tree with the index attributes *did* and *surr* (in this order). Finally *xmlsurrvaldid* denotes the use of the compound B-tree with the index attributes *surr*, *val* and *did*.

We rewrote the query to SQL with the method from section 4.2 and measured both, the elapsed time and the number of pages that were accessed for answering the SQL queries. The number of retrieved pages is further divided into the overall number of physical page accesses, logical

<sup>1</sup><http://www.transaction.de>

accesses to the index pages, and logical accesses to the data pages. The physical page accesses occur whenever the database system requests a page from secondary storage, i.e. the page is not available from database cache. Logical page accesses occur whenever a page is accessed by the database system. One physical page access leads to at least one logical page access. Several logical page accesses occur if the page is accessed several times, e.g., if tuples on pages are repeatedly read in different stages of query processing. In these cases no physical access occurs if the page is available in the database cache.

We analyze selection and projection separately. It would of course be possible to rewrite the query into one large SQL statement, but this would hide the diversity of the query parts and significant facts about the nature of XML queries.

For our example data set (10000 documents) the following numbers of tuples were returned for each processing step (Table 3).

Query Step	Selection (SDID)	Projection (RV)
Number of Tuples	104	554

Table 3: Number of tuples returned for each processing step

## 5.1 Selection

As noted above, we have separated the queries into a selection and projection part. The selection restricts values and compound surrogates and outputs a set of document ids. The set of document ids is then further processed in the projection.

The selection query is graphically illustrated in Figure 4 (due to reasons of visibility the figure only shows 8 point restrictions). The query cuts through the cube as the two dimensions are restricted, while the third is variable. The query is located on two parallel planes orthogonal to the value dimension and is processed using 16 point restrictions (one for each surrogate and value restriction). The compound surrogates are dense in their dimension.

### 5.1.1 UB-Tree *xmlhc*

For the UB-Tree the results of this query are located on the same page with a high propability due to the space-filling Z-curve. As the query is processed by repeatedly stabbing through the threedimensional space pages are read severaly times; caching can be used in this situation to speed up query processing. In this case the caching is intra-query as the query processing itself benefits from the reuse of pages that were already read from secondary storage at earlier stages of query processing.

The consequence is that the average time per page sharply decreases, as the page is processed only in memory. This is shown by the measurements in Figure 4.

Index	log. Idxp.	log. Datp.	phys. Pages	ms/phys. Page	DB size (pages)
UB-Tree	918	782	693	0.3	20428
DidSurr	213	11946	11957	1.7	14426
SurrValDid	39	22	36	0.69	12132

Table 4: Page numbers for the selection

### 5.1.2 Compound B-Tree *xmldidsurr*

For a compound B-Tree on the attributes *did* and *surr* the scenario is completely different compared to the UB-Tree. The query processor of the database system cannot use any restriction on the *did*

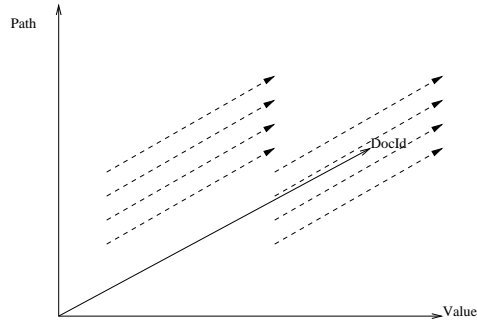


Figure 4: Point restrictions for selection query

attribute, so it has to start with the smallest value for the compound surrogate dimension, starts to read all data pages and performs a post-filtering. Table 4 shows that the query reads almost all pages of the database. The query is slowest among all the others (Table 5). Still it achieves a high page rate. This hints that the query is highly supported from caching. This time it is not intra-query caching as with the UB-Tree, but caching from the operating system.

### 5.1.3 Compound B-Tree *xmlsurrvaldid*

In our third measurement the selection is performed with **xmltriple** indexed with a compound B-Tree on the compound surrogate attribute, the value attribute, and the document id attribute. This index exactly supports the restrictions of the query. The number of physical pages read is only 5 % of physical pages read for the UB-Tree (Table 4).

Index	Selection	Projection
UB-Tree	0.26s	8.73s
DidSurr	20.3s	1.96s
SurrValDid	0.02s	9.18s
Tupel	104 (DocIds)	554 (Tags)

Table 5: Running times for selection and projection queries

## 5.2 Projection

The projection query outputs values and restricts the document ids and the compound surrogates. The restrictions are no range restrictions but point restrictions. The cardinality of the set of points in the document id dimension depends on the result set of the selection query. It is important to note that the result set of the selection query is usually not a range (a range would be quite unlikely). The compound surrogates in the other dimension are restricted to 9 point values (one compound surrogate for title and 8 compound surrogates for keywords). All values which answer the query are consequently located on  $9 \times 104 = 936$  parallel straight lines intersecting the three dimensional space. Figure 5 sketches the scenario with three lines due to the sake of clarity.

### 5.2.1 UB-Tree

Since the straight lines completely stab through the universe (there is no restriction in the value dimension) we can deduce some more information about the processed data by calculating the expected number of page accesses. According to Table 6 the size of the database is 20428 pages. This results in approx.  $2^{15} = 2^{3 \cdot 5}$  pages. This leads to approx.  $2^5 = 32$  pages in each of the three

dimensions meaning that each of the 936 straight lines touches 32 pages. This sums up to 29952 logical data page accesses.

Figure 6 presents the page numbers for the projection query from our performance tests. The accessed logical data pages are not exactly the same as from the calculation above. We assumed a uniform distribution of the data. This is not the case for the data we used.

Another drawback for the UB-Tree in this measurement is the very high number of physical accesses to the pages. The numbers show that for our simple example query almost one third of the database is being read (6441 pages of 20428 pages) although the result set is very small. Most of these physical page accesses are data pages as only 1% of the overall pages in the database are index pages. The high number of physical page accesses is a combination of the way the query is processed (as explained above), the distribution of the selection results which are not ranges but spread over the *docid* dimension and the clustering of the UB-Tree which does not favour one or two attributes but treats all attributes in an equal manner.

Index	log. Idxp.	log. Datp.	phys. Pages	ms/Page	DB size
UB-Tree	39700	31043	6441	1.35	20428
DidSurr	384	436	305	6.42	14426
SurrValDid	6	463	468	19.6	12132

Table 6: Page numbers for the projection

### 5.2.2 Compound B-Tree *didsurr*

The compound index on the attributes *did* and *surr* is the fastest index for projection as the query restricts exactly the index attributes to points. As there are more dids than compound surrogates there are more index pages read than for the *surrvaldid* index.

### 5.2.3 Compound B-Tree *surrvaldid*

Indexing the **xmltriple** relation with a compound B-Tree (index attributes *surr*, *value*, *did*) leads to low page numbers in comparison to the UB-Tree. In contrast the running time of the query is almost the same as for the UB-Tree. This leads to approx. 20 ms per page. A closer look on the way the query is processed reveals that both numbers are reasonable. The projection query restricts on the surrogates and on *did* and there is no restriction on the values. The system starts with reading the data pages starting with the smallest value for *value* and *did* until it terminates when the surrogate range is processed for the keyword surrogates. The system accesses the B-Tree a second time for the title surrogate. The results of the projection query are determined by post filtering the retrieved pages. The retrieved 463 data pages carry approx.  $463 * 100 = 46300$  tuples. The post filtering has to be done for each of the 104 tuples that resulted from the selection. This leads to  $104 * 46300 = 4815200$  overall comparisons. The observation shows that the projection for *surrvaldid* is not I/O bound but CPU bound.

## 6 Summary

We have shown a multidimensional mapping scheme for XML and relational database systems. In addition we have analyzed three different multidimensional indexing methods for our mapping (the UB-Tree and two variants of a compound B-Tree). Due to the special orthogonal requirements for selection and projection currently a combination of the two compound B-Trees seems to be the most promising one.

The use of two indexes requires more maintenance work for insertion as a tuple is effectively inserted two times into the indexes. Another drawback is the consumption of additional storage for the second index.

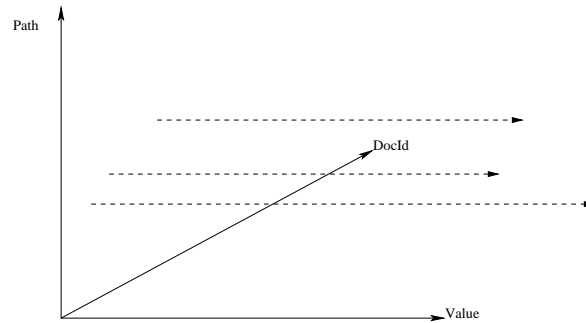


Figure 5: Point restriction for projection query

Our future work on the subject of XML indexing will cover an analysis of secondary indexes and an analysis of the time and space trade-off for answering queries on XML. We will also present a comparison with the edge mapping approach[5].

## References

- [1] Bayer R., *The universal B-Tree for multidimensional Indexing: General Concepts.*, World-Wide Computing and Its Applications '97 (WWCA '97), Tsukuba, Japan, 1997.
- [2] Bayer R., McCreight E., *Organization and Maintenance of Large Ordered Indexes*, Acta Informatica 1, 1972, pp. 173-189.
- [3] Bayer R., *XML Databases: Modelling and Multidimensional Indexing*, DEXA 2001 Invited Talk, Munich, <http://www3.in.tum.de>
- [4] Cooper B., Sample N., Franklin M.J., Hjaltason G.R., Shadmon M., *A Fast Index for Semistructured Data*, VLDB 2001, pp. 341-350.
- [5] Florescu D., Kossmann D., *A Performance Evaluation of Alternative Mapping Schemes for Storing XML Data in a Relational Database*, Rapport de Recherche No. 3680, NRIA, Rocquencourt, France, May 1999.
- [6] Markl V., *MISTRAL: Processing Relational Queries using a Multidimensional Access Technique*, Ph.D. Thesis, TU Munchen, infix Verlag, St. Augustin, DISDBIS 59, ISBN 3-89601-459-5, 1999.
- [7] Markl V., Ramsak F., Bayer R., *Improving OLAP Performance by Multidimensional Hierarchical Clustering*, Proc. of IDEAS Conf., Montreal, Canada, 1999.
- [8] XML Generator, IBM Alphaworks, <http://www.alphaworks.ibm.com>
- [9] DBLP, Uni Trier, <http://dblp.uni-trier.de/>