

# A Method for Online Analytical Processing of Text Data

Akihiro Inokuchi\*

The Institute of Scientific & Industrial Research,  
Osaka University  
8-1 Mihogaoka, Ibaraki, Osaka, 567-0047, Japan  
inokuchi@ar.sanken.osaka-u.ac.jp

Koichi Takeda

Tokyo Research Laboratory, IBM Japan  
1623-14, Shimotsuruma, Yamato, Kanagawa,  
242-8502, Japan  
takedasu@jp.ibm.com

## ABSTRACT

There are increasingly visible demands for structured/ unstructured information integration and advanced analytics. However, conventional database technology has not been able to present a robust and practical implementation of a truly integrated architecture for such purposes. After working on several industrial applications (in particular, in the healthcare and life sciences area), we have identified fundamental issues and technical approaches to tackle the issues. In this paper, we propose data representations and algebraic operations for integrating semantic information (e.g., ontologies) into OLAP systems, which allow us to analyze a huge set of textual documents with their underlying semantic information. The performance of the prototype implementation has been evaluated using real world datasets, and the high scalability and flexibility of our approach have been confirmed with respect to the computation time.

## Categories and Subject Descriptors

H.2 [DATABASE MANAGEMENT]: Database applications

## General Terms

Performance

## Keywords

OLAP, Text Mining

## 1. INTRODUCTION

Since many of business intelligence applications have been incorporating unstructured (primarily textual) information for more context-oriented analysis and decision-making [4],

---

\*This research was conducted when A. Inokuchi was affiliated with Tokyo Research Laboratory, IBM Japan.

database technology has been seriously challenged to ingest, map, store, and access such text-originated information along with the structured information in a way that two types of information can mutually enhance information discovery and analysis capability. One of the most critical problems is that most of semantics underlying the unstructured information (such as ontological hierarchy, synonymous and antonymous relationship) cannot be effectively managed by conventional database systems. Another significant problem is that a rigid schematic representation (and associated queries and analytic processing) of unstructured information often suffers frequent modifications due to updates to dictionary and ontology for adequately categorized words, phrases, and entities described in the unstructured information. Therefore, it is very important to propose a more flexible representation, which reduces the cost and workload of frequent revision and re-population of the database schema.

The multidimensional database technology has been considered for the interactive analysis of large amounts of data for decision making purposes [23, 13, 1, 2, 10, 12]. Multidimensional data models categorize data either as facts with associated numerical measures or as dimensions that characterize the facts. In a retail business, for example, a purchase transaction would be a fact and the purchase amount and price would be measures, and the type of purchased product, the purchase time and location would be dimensions. Queries for Online Analytical Processing (OLAP) aggregate measures over a range of dimensional values to provide results such as the total sales per month of a given product, leading to overall trends. An important feature of the multidimensional data model is to use hierarchical dimensions to provide as much context as possible for the facts. Dimensions are used for selecting and aggregating data at the desired level of detail. Most of traditional multidimensional databases assume that the dimensional hierarchies are balanced and non-ragged trees.

The star and snowflake schemas, which are representative schemas for the multidimensional data model, store data in fact and dimension tables. A fact table holds a row for each fact in the database and it has a column for each measure, containing the measured value for the particular fact, as well as a column for each dimension that contains a foreign key referring to a dimension table for the particular dimension.

When analyzing unstructured information in a multidimensional data model, a document would be typically represented as a fact, and categories of keywords, such as protein, gene, or disease in the life science domain, would be selected as axis for the interactive analysis as shown in Fig-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

CIKM'07, November 6–8, 2007, Lisboa, Portugal.

Copyright 2007 ACM 978-1-59593-803-9/07/0011 ...\$5.00.

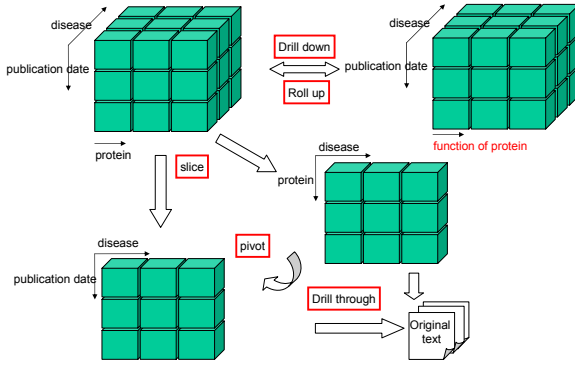


Figure 1: OLAP for Text Data

ure 1. Each cell of the cube in the figure stores the number of the corresponding documents. Operations, such as drill down, roll up, slice, dice, pivoting or drill through, are available for analyzing/aggregating large amounts of documents and their contextual information to obtain insights. It is often very difficult, however, to define a set of dimensions and their hierarchies for a huge set of keywords such as protein name, gene names. For example, the number of distinct keywords in data used in our experiments in Section 5 is 13,640,593. To design a hierarchy for OLAP, we use ontologies such as Unified Medical Language System (UMLS: <http://umlsinfo.nlm.nih.gov>) and the Gene Ontology (GO: <http://www.geneontology.org>), each of which needs to be represented as a kind of a directed acyclic graph, rather than a set of balanced and non-ragged trees. When we assume that each node of the hierarchy corresponds to a dimension, many missing values and a set of multiple values for the node could possibly be introduced. In addition, because the number of nodes in the hierarchy becomes very large and a complex relationship among the nodes exists, we cannot store the data in the star schema and efficiently aggregate the data within the hierarchy under a straightforward implementation.

In this paper, we propose a data representation and algebraic operations to integrate a multidimensional model with ontologies to analyze a huge set of textual documents. This paper describes

- how we design the data representation and its algebraic operations to realize multidimensional model and to integrate with ontologies,
- how we store very "high dimensional data extracted from text documents into a relational database,
- how we efficiently enumerate term and document frequencies for each cell in the cube view, and
- how we can get sufficient performance to provide user interactivity.

The rest of this paper is organized as follows. Section 2 addresses a hierarchy and an ontology. Section 3 defines our data representation and its algebraic operations. In Section 4, we introduce our schemas to efficiently compute the distributions and their implementations. Section 5 presents experiments using about 500,000 abstracts. Section 6 discusses related works. Finally, Section 7 concludes this paper.

## 2. HIERARCHY AND ONTOLOGY

In this section, we give formal definitions of a hierarchy and an ontology according to [3]. If  $S$  is a nonempty set, and  $\leq \subseteq S \times S$ , then  $(S, \leq)$  is an ordering<sup>1</sup>. If  $x \leq x$  for  $x \in S$ , then  $S$  is reflexive. If  $x \leq y$  and  $y \leq z \rightarrow x \leq z$  for  $x, y, z \in S$ , then  $S$  is transitive. If  $x \leq y$  and  $y \leq x \rightarrow x = y$  for  $x, y \in S$ , then  $S$  is anti-symmetric.  $(S, \leq)$  is a partial ordering if  $S$  is a reflexive, transitive, and anti-symmetric binary relation on  $S$ .

**better:** Let  $(S, \leq_1)$  and  $(S, \leq_2)$  be two orderings. We say  $(S, \leq_1)$  is better than  $(S, \leq_2)$  iff  $\forall x, y \in S (x \leq_1 y \rightarrow x \leq_2 y)$ . In addition, we say that  $(S, \leq_1)$  is strictly better than  $(S, \leq_2)$  iff  $(S, \leq_1)$  is better than  $(S, \leq_2)$  and  $(S, \leq_2)$  is not better than  $(S, \leq_1)$ .

**hierarchy:** Let  $(S, \leq)$  be a partial ordering. A hierarchy of  $S$  is an ordering  $(S, \preceq)$  such that (1)  $(S, \preceq)$  is better than  $(S, \leq)$ , (2)  $(S, \preceq)$  is the reflexive, transitive closure of  $(S, \leq)$ , and (3) there is no other ordering  $(S, \sqsubseteq)$  satisfying the above two conditions and that  $(S, \sqsubseteq)$  is strictly better than  $(S, \preceq)$ .

**ontology:** Suppose  $\Sigma$  is some finite set of strings and  $S$  is some set. An ontology w.r.t.  $\Sigma$  is a partial mapping  $\Theta$  from  $\Sigma$  to hierarchies for  $S$ .

**Example 1:** When  $S$  is given as  $\{tire, car, hubcap\}$ , where tire is a part of car, hubcap is a part of car, and hubcap is a part of a tire. In addition, everything is a part of itself. For the set  $S$ , a partial order is defined as  $\{(tire, tire), (car, car), (hubcap, hubcap), (tire, car), (hubcap, car), (hubcap, tire)\}$ , and only one hierarchy is defined as  $\{(tire, car), (hubcap, tire)\}$ .

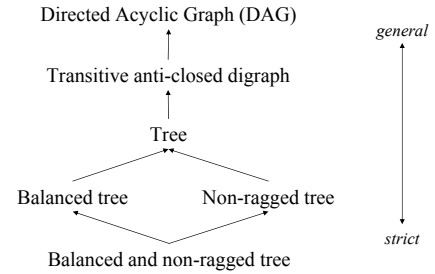


Figure 2: Taxonomy of Hierarchy

Hierarchies can be classified according to their generality as shown in Figure 2 [21].

**DAG:** Directed acyclic graph (DAG) which is a directed graph with no directed cycles is the most general class in the taxonomy of hierarchy. The hierarchy introduced above to define the ontology is a proper subclass of this class.

**transitive anti-closed digraph:** The anti-closure can be produced as follows. If there are an edge of length 1 and a path from node A to node B of length 2 or more, then remove the path of length 1 from A to B. The hierarchy defined above is in this class.

**tree:** A tree is a DAG where each node can only have one parent, except for one node which has no parents and which is called the root.

**balanced tree:** Levels in an unbalanced hierarchy have a consistent parent-child relationship but have a logically inconsistent levels. The hierarchy branches also can have in-

<sup>1</sup>This paper uses  $\leq$  to represent a direct relation between two elements in the set  $S$ ,  $<^+$  to represent its transitive closure, and  $\leq^+$  to represent its transitive closure or represent that the elements are equal.

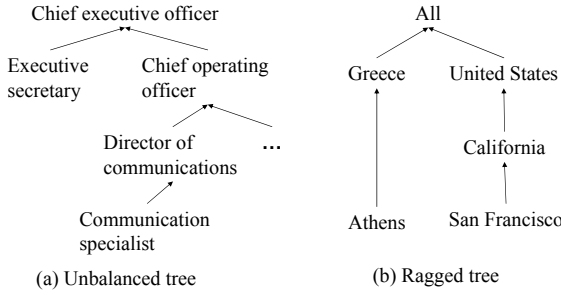


Figure 3: An Unbalanced Tree and a Ragged Tree

consistent depths. For example, an unbalanced hierarchy in Figure 3 (a) shows a chief executive officer on the top level of the hierarchy and at least two of the people that might branch off below including the chief operating officer and the executive secretary. The chief operating officer has more people branching off also, but the executive secretary does not. The parent-child relationships on both branches of the hierarchy are consistent. However, an executive secretary is not the logical equivalent of a chief operating officer<sup>2</sup>. **non-ragged tree:** Each level in a ragged hierarchy has a consistent meaning, but the branches have inconsistent depths because at least one branch has no corresponding node at some level. Figure 3 (b) shows a geographic hierarchy that has Country, State, and City levels defined. The hierarchy becomes ragged when some member does not have an entry at all of the levels. For example, a branch has no entry for the State level because this level is not applicable to Greece for the business model in this example. In this example, the Greece and United States branches have different depths to the leaf, creating a ragged hierarchy. **balanced and non-ragged tree:** Most of traditional multidimensional databases use hierarchies of this class.

### 3. DATA OBJECT AND OPERATIONS

In this section, we give formal definitions of our data representation and operations according to [22].

#### 3.1 Data Object

Given a hierarchy (or an ontology)  $(S, \leq)$ , a fact schema is defined as  $\mathcal{S} = (\mathcal{F}, \mathcal{T})$ , where  $\mathcal{F}$  is a fact type and  $\mathcal{T}$  is a hierarchy type  $\mathcal{T} = (\mathcal{C}, \leq_{\mathcal{T}}, \top_{\mathcal{T}})$  which is strictly better than  $(S, \leq)$  and the relations in  $(S, \leq)$  required for analyzing the documents are remaining in  $\mathcal{T}$ . The hierarchy type is a three-tuple  $(\mathcal{C}, \leq_{\mathcal{T}}, \top_{\mathcal{T}})$ , where  $\mathcal{C} = \{\mathcal{C}_j, j = 1, \dots, n\}$  is a set of category types of  $\mathcal{T}$ , and  $\leq_{\mathcal{T}}$  is a partial order on  $\mathcal{C}$ 's, with  $\top_{\mathcal{T}} \in \mathcal{C}$  being the top element of the ordering. The intuition is that the top element of the ordering logically contains all other elements, that is  $\forall \mathcal{C}_j \in \mathcal{C}, \mathcal{C}_j \leq^+ \top_{\mathcal{T}}$ .

A hierarchy instance  $T$  of type  $\mathcal{T}$  is a two-tuple  $T = (\mathcal{C}, \leq)$ , where  $\mathcal{C}$  is a set of categories  $c_j$  such that  $Type(c_j) = \mathcal{C}_j$ , and  $\leq$  is a partial order on  $\mathcal{C}$ . Functions to give the set of immediate predecessors and successors of a category  $c_j$  are defined as  $pred : \mathcal{C} \rightarrow 2^{\mathcal{C}}$  and  $succ : \mathcal{C} \rightarrow 2^{\mathcal{C}}$ . That is,  $pred(c_j) = \{c' \mid c' > c_j\}$  and  $succ(c_j) = \{c' \mid c' < c_j\}$ , respectively. Each category  $c \in \mathcal{C}$  has an associated set  $dom(c)$  called its domain. The members of  $dom(c)$  are called

values of the category  $c$ . An element in  $dom(c)$  is represented as  $c : v$ . In addition, a function  $below$  to give a set of values is also defined as  $below(c) = \{dom(c') \mid c' \leq^+ c\}$ .

**Example 2:** We have a hierarchy instance  $T$ , a part of which is depicted in Figure 4 (a). Categories such as “Software”, “OS”, “Middleware”, “Application”, “Windows”, “Linux”, “AIX” are contained in  $\mathcal{C}$ .  $pred(OS)$  has only one element “Software”, and  $succ(OS)$  contains {“Windows”, “Linux”, “AIX”}.  $dom(Windows)$  contains “Windows XP”, “Windows Me”, “Windows 2000”, and so on.  $below(All)$  contains all values of all categories.

Let  $F = \{f_i, i = 1, \dots, m\}$  be a set of facts. A fact-hierarchy relationship between  $F$  and  $T$  is a set  $R = \{(f, c : v)\}$ , where  $f \in F$ ,  $c \in \mathcal{C}$ , and  $v \in dom(c)$ . Thus,  $R$  links facts to hierarchical values. We say that fact  $f$  is characterized by a hierarchical value  $c : v$ , written by  $f \rightsquigarrow c : v$ , if  $\exists c' \in \mathcal{C} ((f, c' : v') \in R \wedge c' \leq^+ c \wedge v = v')$ . Our data object is a four tuple  $D = \{S, F, T, R\}$ , where  $S = (\mathcal{F}, \mathcal{T})$  is the fact schema,  $F$  is a set of facts where  $Type(f) = \mathcal{F}$ ,  $T = (\mathcal{C}, \leq)$  is a hierarchy instance where  $Type(c_j) = \mathcal{C}_j$  for  $c_j \in \mathcal{C}$  and  $\mathcal{C}_j \in \mathcal{C}$ , and  $R$  is a set of fact-hierarchy relations such that  $(f, c : v) \in R \Rightarrow f \in F \wedge \exists c \in \mathcal{C} (v \in dom(c))$ .

**Example 3:** We have the hierarchy instance  $T$  and an analyzed document which is depicted in Figure 4 (b).  $F$  contains a set of document identifiers. Terms in the document whose document ID is 1 in Figure 4 are annotated in preprocessing, e.g., categories “windows” and “workstation” are assigned to terms “windows 2000” and “IntelliStation 6217”, respectively, and  $(1, windows : windows\_2000)$  and  $(1, workstation : IntelliStation\_6217)$  are stored in  $R$ .

Conceptually,  $R$  corresponds to a relation  $R' \subseteq 2^{dom(c_1)} \times \dots \times 2^{dom(c_n)}$  which is not a normalized relation.  $R'$  corresponds to a fact table for a star schema, and each row and column in  $R'$  correspond to a document (fact) and a category (dimension value in the star schema), respectively. A naive method cannot store the data in a relational database and efficiently aggregate the data along the hierarchy, because the relation has many missing values and a set of multiple values for each attribute  $c_j$ , the number of attributes in the relation becomes very large and a complex relationship among the attributes (columns) exists.

**Example 4:** A hierarchy instance  $T$  used in Section 5 is a transitive anti-closed digraph which has more than 240,000 nodes (categories) and more than 340,000 edges and whose depth is 24. If the hierarchy instance is a tree,  $V = E + 1$ , where  $V$  and  $E$  are the numbers of categories and edges of the hierarchy, respectively. However, because the difference between the numbers of the categories and edges in  $T$  is so large, the hierarchy instance  $T$  used in Section 5 have a very complex relationship. In addition, about 36,400,000 elements in a conceptual relation  $R'$  have values. Since the number of attributes,  $n$ , for  $R'$  is greater than 240,000 and  $R'$  has more tuples than 500,000 in Section 5, most of elements in  $R'$  are missing values. Furthermore, more than 7,600,000 elements in  $R'$  has a set of values, and  $dom(T)$  has about 193,000,000 distinct values.

The function  $g(c)$  is defined as a user-defined function to return a set of fact-hierarchy relations, and another function  $G(g(c))$  is defined as  $G(g(c)) = \{(c : v) \mid (f, c : v) \in g(c)\}$ . For  $k = 1, \dots, q$  and  $(c_k : v_k) \in G(g_k(c_k))$ , we define a function  $Group$  as  $Group(c_1 : v_1, \dots, c_q : v_q) = \{f \mid f \in F \wedge \bigwedge_{k=1, \dots, q} (f, c_k : v_k) \in g_k(c_k)\}$ . These functions are used to aggregate the distributions of documents for each

<sup>2</sup>[http://publib.boulder.ibm.com/infocenter/db2luw/v8/index.jsp?topic=/com.ibm.db2.udb.db2\\_olap.doc/cmdhierarchy.htm](http://publib.boulder.ibm.com/infocenter/db2luw/v8/index.jsp?topic=/com.ibm.db2.udb.db2_olap.doc/cmdhierarchy.htm)

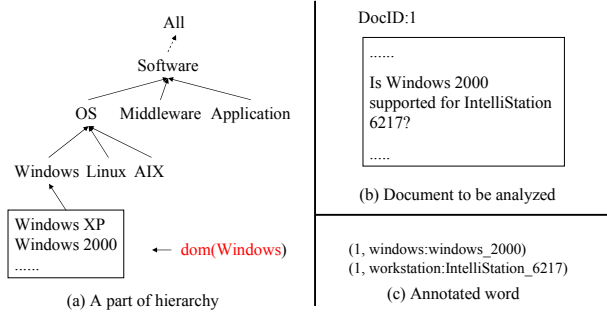


Figure 4: Examples of a Hierarchy and Fact-Hierarchy Relations

$$\begin{aligned}
 g^{(1)}(c) &= \{(f, c' : v') \mid (f, c' : v') \in R \wedge c = c' \wedge v' = \text{dom}(c')\} \\
 g^{(2)}(c) &= \{(f, c : c'') \mid (f, c' : v') \in R \wedge c'' \in \text{succ}(c) \\
 &\quad \wedge c' \leq^+ c'' \wedge v' \in \text{dom}(c')\} \\
 g^{(3)}(c) &= \{(f, c' : v') \mid (f, c' : v') \in R \wedge c' \leq^+ c \wedge v' \in \text{dom}(c')\}
 \end{aligned}$$

Figure 5: Examples of User-Defined Functions

keyword and category. For example,  $g(c)$  is provided as functions  $g^{(1)}, g^{(2)}, g^{(3)}$  ( $g \in \{g^{(1)}, g^{(2)}, g^{(3)}, \dots\}$ ) as shown in Figure 5. First,  $g^{(1)}(c)$  is used to aggregate the distributions for keywords belonging to the specified category  $c$ . The second function  $g^{(2)}(c)$  is used to aggregate the distributions for the immediate successors (subcategories) of the specified category  $c$ . The third one  $g^{(3)}(c)$  is used to aggregate the distributions for keywords belonging to *below*( $c$ ). Users can define any additional functions as required for the intended analysis of a set of documents.

### 3.2 Operations

This subsection defines operations on our data object.

**selection  $\sigma'$ :** Given a compound predicate  $P = p_1 \text{ or } \dots \text{ or } p_l$  where each atomic predicate  $p_i$  is represented in the form of  $c : v \text{ or } c : *$ . The selection  $\sigma'$  is defined as  $\sigma'_P(D) = (S, F', T, R')$ , where  $F' = \{f \mid f \in F \wedge (f \rightsquigarrow p_1 \vee \dots \vee f \rightsquigarrow p_l)\}$ , and  $R' = \{(f', c : v) \in R \mid f' \in F'\}$ . For example, a set of documents having any keywords belonging to a category 'software' is given by  $\sigma'_{\text{software}':*}(D)$ . Other examples  $\sigma'_{\text{gene\_name}':\text{'BIKE'}}(D)$  and  $\sigma'_{c:v_1}(\sigma'_{c:v_2}(D))$  represent a set of documents having a term 'BIKE' belonging to a category 'gene\_name' and a set of documents having a term  $v_1$  and  $v_2$  belonging to a category  $c$ , respectively.

**difference  $-$ :** Given two data objects  $D_1 = (S_1, F_1, T_1, R_1)$  and  $D_2 = (S_2, F_2, T_2, R_2)$  such that  $S_1 = S_2 = S$ , the difference is defined as  $(S, F_1, T_1, R_1) - (S, F_2, T_2, R_2) = (S, F', T_1, R')$ , where  $F' = F_1 - F_2$  and  $R' = \{(f, c : v) \mid f' \in F', (f', c : v) \in R\}$ . For example, a set of documents which has a term  $v_1$  and does not have the term  $v_2$  is  $\sigma'_{\top:v_1}(D) - \sigma'_{\top:v_2}(D)$ .

**projection  $\pi'$ :** The projection is defined as  $\pi'_{c_1 \vee \dots \vee c_l}(D) = (S, F, T, R')$ , where  $R' = \{(f, c : v) \in R \mid f \in F \wedge (f \rightsquigarrow c_1 : * \vee \dots \vee f \rightsquigarrow c_l : *)\}$ .

**aggregation  $\alpha$ :** Given a set of categories and functions  $T = (c_1, \dots, c_q, g_1, \dots, g_q)$ , the aggregation  $\alpha$  is defined as  $\alpha[T, 'count'](D) = (S', F', T', R')$ , where  $S', F', T'$ , and  $R'$  are defined in Figure 6.  $('count', \top)$  represents a partial relation  $'count' < \top$ .

**Example 5:** We give some examples of the aggregation

$$\begin{aligned}
 S' &= (F', T') \\
 F' &= 2^F \\
 T' &= (C', \leq'_{T'}, \top_{T'}) \\
 C' &= C \cup \{'count'\} \\
 \leq'_{T'} &= \leq_T \cup \{('count', \top)\} \\
 \top_{T'} &= \top_T \\
 F' &= \{Group(c_1 : v_1, \dots, c_q : v_q) \mid (c_1 : v_1, \dots, c_q : v_q) \\
 &\quad \in G(g_1(c_1)) \times \dots \times G(g_q(c_q)) \\
 &\quad \wedge Group(c_1 : v_1, \dots, c_q : v_q) \neq \emptyset\} \\
 T' &= (C', \leq') \\
 C' &= C' \cup \{'count'\} \\
 \leq' &= \leq_T \cup \{('count', \top)\} \\
 R' &= R'_1 \cup R'_2 \\
 R'_1 &= \{(f', c' : v') \mid \exists (c_1 : v_1, \dots, c_q : v_q) \in G(g_1(c_1)) \times \dots \\
 &\quad \times G(g_q(c_q)) (f' = Group(c_1 : v_1, \dots, c_q : v_q) \\
 &\quad \wedge f' \in F' \wedge \exists k (c_k : v_k = c' : v'))\} \\
 R'_2 &= \bigcup_{(c_1 : v_1, \dots, c_q : v_q) \in G(g_1(c_1)) \times \dots \times G(g_q(c_q))} \{(s, 'count' : |s|) \mid \\
 &\quad s = Group(c_1 : v_1, \dots, c_q : v_q) \wedge s \neq \emptyset\}
 \end{aligned}$$

Figure 6: Definition of Aggregation

$\alpha$ , and will give detailed procedures using a Document-Term Matrix in Section 5.2. Table 1 shows a top N ranking of protein names based on document frequencies. It is obtained by calculating  $\alpha[\text{'protein'}, g, 'count1'](\sigma'_{\text{'disease'}':\text{'diabetes'}}(D))$ , sorting the result in descending order of frequency, and fetching N rows. Table 2 shows the results of the query  $\alpha[\text{'company'}, \text{'GeneSymbol'}, g_1, g_2, 'count2'](D)$ . If the result is analyzed for a set of patent documents, a strategist for a pharmaceutical company might be able to find associations between companies and genes. The query  $\alpha[\text{'protein'}, \text{'protein'}, g_1, g_2, 'count3'](D)$  may be useful to find interactions between proteins. Note that we can select the same categories as cube's axes unlike the traditional multidimensional database.

Table 1: Top N ranking

| protein                   | # of documents |
|---------------------------|----------------|
| Flavohemoprotein          | 347            |
| Lamin L                   | 240            |
| Insulin                   | 151            |
| nterferon gamma precursor | 97             |
| ...                       | ...            |

Table 2: 2 Dimensional Map

|          | LEPR     | TM4SF2   | INS      | ADAM2    | ... |
|----------|----------|----------|----------|----------|-----|
| company1 | $x_{11}$ | $x_{12}$ | $x_{13}$ | $x_{14}$ | ... |
| company2 | $x_{21}$ | $x_{22}$ | $x_{23}$ | $x_{24}$ | ... |
| company3 | $x_{31}$ | $x_{32}$ | $x_{33}$ | $x_{34}$ | ... |
| ...      | ...      | ...      | ...      | ...      | ... |

Other operations are similarly defined, although they are omitted because of lack of space.

By using above operators, we will show how common OLAP operators can be defined.

**roll-up & drill-down:** In the traditional multidimensional database, there are two types of rolling up operation, one is dimensional rolling up and the other is hierarchical rolling up. For example, let  $S$  be the fact table  $S(\text{product}, \text{city}, \text{time})$ ,

sale), where sale is a measure, and  $L(city, state, country)$  be one of the dimension tables. The dimensional rolling up is represented as  $product, city \chi_{product, city, SUM(sale)}(S)$  in the case of one dimension being dropped, and  $city \chi_{city, SUM(sale)}(S)$  in the case of two dimensions being dropped<sup>3</sup>. The hierarchical rolling up is represented as

$product, state, time \chi_{product, state, time, SUM(sale)}(S \bowtie_{city} L)$ . It is possible to define more than  $2^k$  roll-up queries for the  $k$  dimensions of the traditional multidimensional database. In our case, the dimensional rolling up corresponds to moving from  $\alpha[c_1, \dots, c_q, g_1, \dots, g_q, 'count'](D)$  into  $\alpha[c_1, \dots, c_h, c_{h+2}, \dots, c_q, g_1, \dots, g_h, g_{h+2}, \dots, g_q, 'count'](D)$  for the case of one dimension being dropped, and the hierarchical rolling up corresponds to moving from  $\alpha[c_1, \dots, c_q, g_1^{(1)}, \dots, g_q^{(1)}, 'count'](D)$  to  $\alpha[c_1, \dots, c_h, c'_{h+1}, c_{h+2}, \dots, c_q, g_1^{(1)}, \dots, g_h^{(1)}, g_{h+1}^{(2)}, g_{h+2}^{(1)}, \dots, g_q^{(1)}, 'count'](D)$ , where  $c'_{h+1} \in pred(c_{h+1})$ .

**slice & dice:** The slice operation performs a selection on one dimension of the given cube, resulting in a subcube, and the dice operation defines a subcube by performing a selection on two or more dimension. For example, in the traditional multidimensional database, slice and dice operations are represented as  $city, time \chi_{city, time, SUM(sale)}(\sigma_{product=p_1}(S))$  and  $product, city, time \chi_{product, city, time, SUM(sale)}(\sigma_P(S))$ , where  $P = (product \in \{p_1, p_2\} \text{ and } city \in \{c_3, c_4\})$ . In our case, the slice is represented as  $\alpha[T, 'count'](\sigma_p(D))$ , and the dice is represented as  $\alpha[T, 'count'](\sigma'_{p_1 \text{ or } p_2}(\sigma'_{p_3 \text{ or } p_4}(D)))$ .

**pivot:** The pivot operation is a visualization operation that rotates the data axes in view in order to provide an alternative presentation of the data, which corresponds to moving from  $\alpha[c_1, c_2, g_1, g_2, 'count']$  into  $\alpha[c_2, c_1, g_2, g_1, 'count']$ .

## 4. IMPLEMENTATION

A key strategy for speeding up cube view processing for the traditional multidimensional database is to use pre-computed cube views. The pre-computation reduces the response times to queries potentially involving huge amounts of data and allows interactive data analysis in the traditional approaches. However, it is impossible to pre-compute or pre-aggregate in advance of receiving queries for all of the combinations of values in our situation, because the situation where each document has many values and there are a lot of categories is combinatorially explosive. For example, the average number of annotated terms which each documents have is about 380 and the number of categories is more than 240,000 for the data used in Section 5.

In this section, we design table schemas and data structures to achieve query response times that are as fast as possible. Since a hierarchy for analyzed documents constitutes a transitive anti-closed digraph rather than a set of balanced and non-ragged trees, it cannot be stored in a star schema or snowflake schema. For computation efficiency in aggregating the distributions of documents, the hierarchy is indexed as follows. A depth first search traverses the hierarchy from root category  $c_{root}$  whose type  $Type(c_{root})$  is equal to  $\tau_T$  assigning a preorder, postorder, and depth to each category, and it backtracks if and only if it reaches leaf nodes. This means that it does not backtrack when it reaches any internal nodes which it has already visited.

The assigned preorders and postorders make it possible to

<sup>3</sup>  $a \chi_{a, sum(b)}$  represents an SQL query “SELECT a, SUM(b) FROM ... GROUP BY a”.

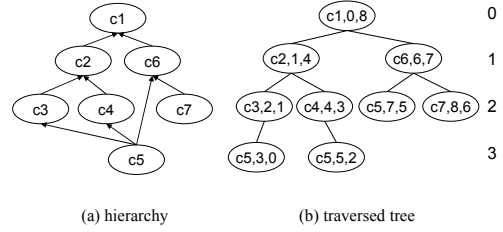


Figure 7: A hierarchy and a Traversed Tree

handle ancestor-descendant containment in the hierarchy [7]. In other words, it can check the containment by assigning a preorder and a postorder to each node in a hierarchy and comparing the preorder and postorder values assigned to given two nodes. If a node A is an ancestor of a node B,

$$A's \text{ preorder} < B's \text{ preorder} \ \& \ A's \text{ postorder} > B's \text{ postorder}. \quad (1)$$

We use the preorder-postorder method to index the category hierarchy, since we assume that updates to a category hierarchy (ontology) happens much less frequently than document insertions that the cost of preorder and postorder recalculations is negligible.

**Example 6:** The hierarchy in Figure 7 (a) is traversed to return a tree shown in Figure 7 (b) where each node has a category, an assigned preorder, postorder, and depth. In this figure, all descendants of a category  $c_2$  have preorders which are greater than the preorder of  $c_2$  and have postorders which are less than the postorder of  $c_2$ . We call the tree in Figure 7 (b) a traversed tree.

We define two tables, **CATEGORY H** and **KEYWORD V**, to store the traversed tree and fact-hierarchy relations as

| CATEGORY (CATEGORYNAME CHARACTER, |               |
|-----------------------------------|---------------|
| PATH                              | CHARACTER,    |
| PREORDER1                         | INTEGER,      |
| PREORDER2                         | INTEGER,      |
| PARENT                            | INTEGER), and |
| KEYWORD (ID                       |               |
| PREORDER                          | INTEGER,      |
| VALUE                             | CHARACTER),   |

respectively. Each record in the table  $H$  corresponds to a node in the traversed tree and **CATEGORYNAME**, **PATH**, **PREORDER1**, **PREORDER2**, and **PARENT** in  $H$  are a name of the category, a path from the root node to the corresponding node, a preorder, a value for a postorder plus a depth, and a preorder of its parent of the corresponding node. The reason why we use **PREORDER2** instead of the postorder is that we can check ancestor-descendant containment in the hierarchy as

$$A's \text{ preorder1} (= A's \text{ preorder}) < B's \text{ preorder} \leq A's \text{ preorder2} = A's \text{ postorder} + A's \text{ depth} \quad (2)$$

instead of using condition (1). Each record in the table  $V$  corresponds to  $(f, c : v)$  in  $R$ , and **ID**, **PREORDER**, and **VALUE** in the table  $V$  are a document ID  $f$ , a preorder of the category  $c$ , and a value  $v$  in  $dom(c)$ , respectively.

By using these tables, we can implement the operations introduced in Section 3.2. In the following definitions,  $c$  is provided as input. Although there are multiple records whose values of **CATEGORYNAME** in  $H$  are  $c$ , a record arbitrarily chosen from the records is used in the following operations. In other words, “ $\sigma_P(H)$ ” for  $P = (categoryname = c)$  is



$$\begin{aligned}
g^{(1)}(c) &= \pi_{id, preorder, value}(V \bowtie_{preorder1=preorder} H_c) \\
g^{(2)}(c) &= \pi_{id, H.preorder1, H.catename \text{ as } value} ( \\
&\quad (H_c \bowtie_{H.parent=H_c.preorder1} H) \\
&\quad \bowtie_{H.preorder1 \leq V.preorder \leq H.preorder2} V) \\
g^{(3)}(c) &= \pi_{id, preorder, value}(V \bowtie_{preorder1 \leq preorder \leq preorder2} H_c)
\end{aligned}$$

Figure 8: Implementation of User-defined Functions

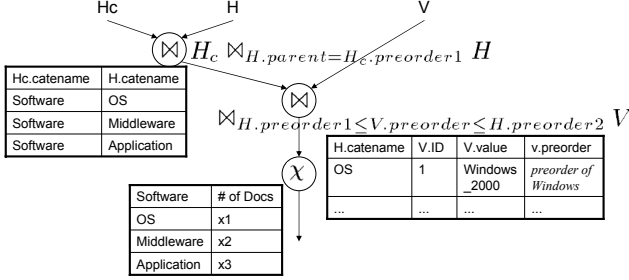


Figure 9: An Example of Aggregation

replaced into “ $\sigma_P(H)$  FETCH FIRST 1 ROWS ONLY”, denoted as  $H_c$ . The choice has no influence on its result.

**selection  $\sigma'$ :** The selection is defined as  $\sigma'_{c:v}(D) = (H, V')$ , where  $V' = \sigma_{id \text{ in } I(V)}$ ,  $I = \pi_{id}(V \bowtie_P H_c)$  and  $P = (preorder1 \leq preorder \leq preorder2 \text{ and } value = v)$ . In the definition, the condition (2) is used as “ $preorder1 \leq preorder \leq preorder2$ ”. Although this calculation needs to join  $H$  with  $V$ , this calculation runs as fast as the selection of  $V$ , since only one record is returned from  $H_c$ .

**difference  $-$ :** The difference is defined as  $(H, V_1) - (H, V_2) = (H, V')$ , where  $V' = \sigma_{id \text{ in } (\pi_{id}(V_1) - \pi_{id}(V_2))}(V_1)$ .

**projection  $\pi'$ :** The projection is defined as  $\pi'_c(D) = (H, V')$ , where  $V' = \pi_{id, preorder, value}(V \bowtie_P H_c)$  and  $P = (preorder1 \leq preorder \leq preorder2)$ .

**aggregation  $\alpha$ :** The aggregation is defined as

$\alpha[(c_1, \dots, c_q, g_1, \dots, g_q), 'count'](D) = value_1, \dots, value_q \chi_{value_1, \dots, value_q, count(distinct id)}(X)$ , where  $X = g_1(c_1) \bowtie_{id=id} \dots \bowtie_{id=id} g_q(c_q)$ . The user-defined functions  $g^{(1)}(c)$ ,  $g^{(2)}(c)$ , and  $g^{(3)}(c)$  are defined in Figure 8.

**Example 7:** Figure 9 shows how to compute  $\alpha[Software', g^{(2)}, 'count'](D)$  for the data shown in Figure 4. First,  $H_c \bowtie_{H.parent=H_c.preorder1} H$  returns combinations  $A_1$  of the category “Software” and its child categories. Second,  $A_1 \bowtie_{H.preorder1 \leq V.preorder \leq H.preorder2} V$  outputs combinations  $A_2$  of the child categories and the corresponding document IDs. The output contains a tuple in the first row, because “Windows\_2000” is a descendant of “OS” in the category hierarchy in Figure 4 and preorders for “Windows\_2000” and “OS” fulfill the condition (2). Finally, the distribution  $H.categoryname \chi_{H.categoryname, count(distinct V.id)}(A_2)$  for the immediate subcategories of “Software” is returned.

## 5. EXPERIMENTS

### 5.1 Data and Preprocess

For testing, we used biomedical documents from MEDLINE ([http://www.nlm.nih.gov/databases/databases\\_medline.html](http://www.nlm.nih.gov/databases/databases_medline.html)). Life science researchers typically use MEDLINE, a bibliography database that covers the biomedical area. MEDLINE is administered by the National Center for Biotechnology Information (NCBI: <http://www.ncbi.nlm.nih.gov>) of the United States National Library of Medicine (NLM: <http://www.nlm.nih.gov/>). It contains approximately 17

million biomedical citations, dating from the mid-1960s up to the present time. Citations in MEDLINE are collected from over 5,000 biomedical journals published worldwide. Biomedical citations in MEDLINE are available to the general public at the PubMed (<http://www.ncbi.nlm.nih.gov/entrez/>). We selected 503,989 abstracts from Medline which contain structured information such as authors and Mesh Terms and unstructured information such as titles and abstracts.

To prepare a fact-hierarchy relation from the documents, the documents written in English are parsed by CCAT [5], a shallow syntactic parser. Because this is a general-purpose parser that has not been trained for biomedical documents, it is difficult to obtain optimized results by parsing documents from various domains [27]. We solve this problem by first annotating the text with domain dictionaries. The annotations facilitate the parsing of the documents even when the parser has not been specifically trained for the domains.

In the first step of the preprocessing, the term annotator finds words in the input text using the term dictionary and identifies these words with their canonical forms. The term dictionary contains pairs of surface forms and canonical forms. For example, most of the technical terms in the medical domain are compound words. The compound noun “repetitive sequence-based polymerase chain reaction” consists of an adjective (repetitive), a past participle of a verb (sequence-based) and three nouns (polymerase, chain, reaction). Thus, biomedical terms tend to consist of a combination of numerals, symbols, and verbs, making it very difficult to find term boundaries. In addition, there are often considerable numbers of expressions that are synonymous with a particular technical term. These can arise from abbreviations or acronyms as well as from spelling variations. If these variations are recognized as different entities, it can often cause problems when aggregating unstructured information in documents. For instance, “DNA” and “deoxyribonucleic acid” are synonyms. The dictionary contains spelling and abbreviation variants and their canonical forms. By reducing these variants to a single canonical form, we treat them as the same entity.

In the second step, the text annotated with a technical term dictionary is passed to the syntactic parser. The parser outputs segments of phrases labeled with their syntactic roles, for example NP (noun phrase) or VG (verb group). In the third step, the category annotator assigns categories to the terms in these segments and phrases. The category dictionary consists of a set of canonical forms and their categories, which also indicate the node labels in the category hierarchy (ontology). A category assigned to each term is an internal node or leaf in the hierarchy.

**Example 8:** Figure 10 shows an example of the preprocessing of a sentence. When “Repetitive sequence-based polymerase chain reaction effects deoxyribonucleic acids” is given as input, an annotator assigns “DNA” as canonical and “proper noun” as part-of-speech to “deoxyribonucleic acids”. After parsing the annotated text, categories are assigned to each term. In Figure 10, “.A.1.2.23.4” represents PATH from a root node to the corresponding node in a category hierarchy.

After preprocessing the 503,989 abstracts, the numbers of  $(f, c : v)$ , records in the table  $H$ , and distinct canonical forms of terms were 193185919, 340154, and 13640593, respectively. The categories contain categories for publication dates, authors, affiliations for the authors, and so on.

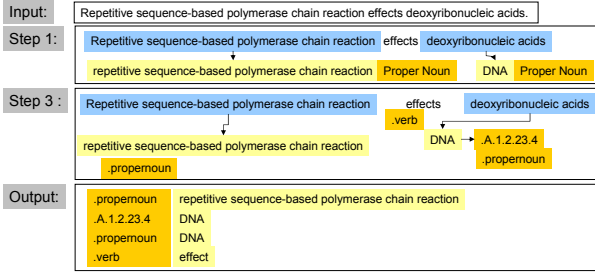


Figure 10: Example of the Preprocessing

## 5.2 Implementation using a DTM

To compare with the method mentioned in Section 4, we used a method with a Document-Term Matrix (DTM) as a proprietary method. In general, a proprietary algorithm and index can compute faster than a method with a persistent store, *e.g.*, the method in Section 4, although it is more difficult to add some functions into the proprietary method and to integrate the proprietary method with other systems compared to the persistent-store method. We explain the method using the proprietary algorithm and index, and the next subsection will explain that the method using the persistent store is comparable to the proprietary method.

We focus how to compute  $\alpha[c, g, count](\sigma'_P(D))$  by a DTM using a simple example, because this is the most fundamental computation. Let the sets of terms and documents be  $T = \{t_1, t_2, \dots, t_n\}$  and  $D = \{d_1, d_2, \dots, d_m\}$ , respectively. A DTM is a matrix  $M = (m_{ij})$  of  $m \times n$ , and an elements  $m_{ij}$  represents how many times the term  $t_j$  appears in the document  $d_i$ . Although storing the whole matrix requires a lot of memory, it can be compressed by storing a pair for each element that is not zero and its corresponding index in each row or column, because the matrix is very sparse.

As mentioned in the previous section, since some categories are assigned to each term, we use a modified DTM. Rows in our DTM correspond to a set of documents  $D = \{d_1, \dots, d_m\}$  similar to the conventional DTM, and columns correspond to a set of pairs of categories and terms  $P = \{c_j : v_{jk} \mid v_{jk} \in \text{dom}(c_j), j = 1, \dots, n\} \cup \{c_j : * \mid j = 1, \dots, n\}$ . The value  $c_j : *$  is used to facilitate aggregating the number of documents for each subcategory, and an element for  $d_i$  and  $c_j : *$  is not zero when  $d_i \sim c_j : *$ .

Figure 11 presents how the method using DTM computes the results for  $\alpha[c_3, g^{(1)}, count](\sigma'_{c_1:v_{11}}(D))$ , when a user has specified the category  $c_3$  after narrowing down to the documents containing the term  $v_{11}$  whose category is  $c_1$ . First, the method narrows the search down to the document set  $\{d_2, d_6, d_{10}\}$  containing  $(c_1 : v_{11})$  (1). In parallel with this Process 1, a set of terms  $\{c_3 : v_{31}, c_3 : v_{32}, c_3 : v_{33}, c_3 : v_{34}, c_3 : v_{35}\}$  whose category is  $c_3$  is output (2). After Processes 1 and 2, the distribution of the documents for the terms appearing in the documents  $\{d_2, d_6, d_{10}\}$  is returned as  $\{(c_3 : v_{31}) : 2, (c_3 : v_{35}) : 1\}$  (3). Process 3 requires much more computation time than Processes 1 and 2. For example, when the user selects a “common\_noun” category, Process 2 returns 340,154 terms for the dataset used in the experiments described in Section 5.3.

When the user specified a category which is an internal node in the category tree, it also computes the distribution of the documents for each subcategory of the specified category, which corresponds to  $\alpha[c_3, g^{(2)}, count](\sigma'_P(D))$ . In

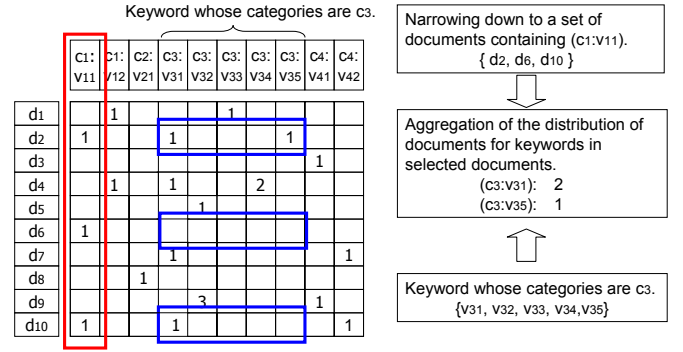


Figure 11: Implementation using a DTM

this case, a set of categories  $\{c'_3 : * \mid c'_3 \in \text{succ}(c_3)\}$  is returned in Process 2.

## 5.3 Experimental Comparison

The method in Section 4 was implemented in Java. It generates SQL queries and accesses a relational database via JDBC (Java Database Connectivity). The method in Section 5.2 was implemented in C++ to compare with the above method. For the evaluation, an IBM IntelliStation with Windows XP, an Opteron-2.2 GHz CPU, and 2 GB of main memory was used. The efficiency of our approach has been confirmed with respect to the computation time.

Figure 12 shows the results of the response time for a query  $\alpha[c, g, count](D)$  for all the 340,154 categories. The DB and DTM in the figure correspond to implementations of the methods mentioned in Section 4 and Section 5.2, respectively. KW and SUB represent the cases where  $g^{(1)}$  and  $g^{(2)}$  as  $g$  in  $\alpha[c, g, count](D)$  were used, respectively. “DTM SUB 1k” shows the results for 1,000 documents sampled randomly from all of the documents. The method of the DB does not compute for the sampled documents but for all of the documents. Each point  $(x, y)$  in the figure means that the method returns the result within  $x$  seconds for  $y\%$  of all of the 340,154 categories. The ideal method is at the upper left corner. While DB could return the result for about 89% of the categories within 0.1 second, DTM could only return results for about 60% of the categories for 1,000 sampled documents, and for about 0.01% of categories for 10,000 sampled documents for KW. In addition,  $\alpha[c, g, count](\sigma_{T: \text{cancer}}(D))$  was calculated for the various categories, and the results for 11,914 documents containing “cancer” were similar to Figure 12. As shown in Figure 12, DB is superior to DTM for most of the categories.

From another viewpoint of the empirical 8-second rule saying that a webpage should be loaded within 8 seconds of a request [20], we can conclude that the better method is the one whose coverage rate in about 10 seconds is better. Figure 13 shows the coverage rates of DB and DTM - between 99.97% and 100% in 5-10 second response time for the query  $\alpha[c, g, count](\sigma_{T: \text{cancer}}(D))$ . Figure 13 shows that DB is inferior to DTM within 8-second constraint. Tables 3 and 4 summarize the experimental results for  $\alpha[c, g, count](D)$  and  $\alpha[c, g, count](\sigma_{T: \text{cancer}}(D))$ , respectively. Although the average computation time of DB is lower than for DTM, the number of categories for which DB cannot return within 10 seconds may become greater than for DTM.

As shown in Table 4 and Figure 13, the averages of the

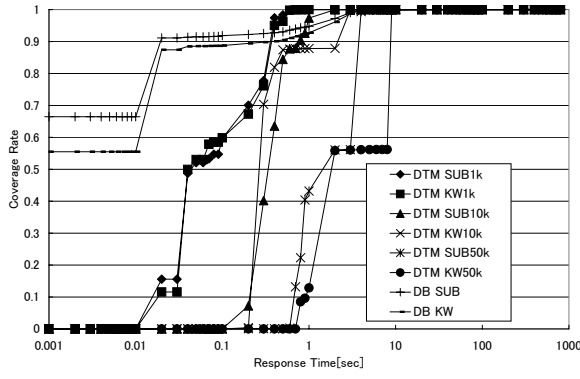
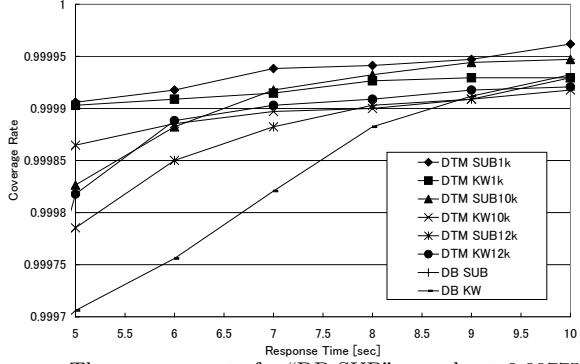


Figure 12: Results for All Documents



\*The coverage rate for "DB SUB" was about 0.99775.

Figure 13: Experimental Results for Documents Containing "cancer"

computation times for DB are superior to DTM. However, the number of categories for which DB cannot respond within 10 seconds is greater than for DTM. When documents containing "cancer" are selected, an SQL query to obtain  $V'$  for  $\sigma_{T:'cancer'}(D) = (H, V')$  is represented as

$$V' = \sigma_{id \text{ in } I(V)} = \pi_{id, V^{(b)}.preorder, V^{(b)}.value}(V^{(a)} \bowtie_{P_2} V^{(b)}), \quad (3)$$

where  $I = \pi_{id}(V \bowtie_{P_1} H_c)$ ,  $P_1 = (value = 'cancer')$ ,  $P_2 = (V^{(a)}.id = V^{(b)}.id \text{ and } V^{(a)}.value = 'cancer')$ , and  $V = V^{(a)} = V^{(b)}$ . The query (3) is represented as

```
SELECT Vb.ID, Vb.PREORDER, Vb.VALUE
FROM V AS Va, V AS Vb
WHERE Va.ID=Vb.ID AND Va.VALUE='cancer'.
```

Table 3: Comp. Times for All Documents

| Method      | avg. comp. time | #10 | #100 |
|-------------|-----------------|-----|------|
| DTM KW 1k   | 0.143 sec.      | 1   | 0    |
| DTM KW 5k   | 0.284 sec.      | 26  | 0    |
| DTM KW 10k  | 0.526 sec.      | 16  | 0    |
| DTM KW 50k  | 4.293 sec.      | 75  | 0    |
| DB KW       | 0.185 sec.      | 10  | 0    |
| DTM SUB 1k  | 0.138 sec.      | 2   | 1    |
| DTM SUB 5k  | 0.176 sec.      | 2   | 0    |
| DTM SUB 10k | 0.405 sec.      | 5   | 1    |
| DTM SUB 50k | 2.091 sec.      | 33  | 0    |
| DB SUB      | 0.137 sec.      | 20  | 2    |

#10 means the number of categories for which the results are not returned within 10 seconds.

Table 4: Computation Times for Documents Containing "cancer"

| Method      | avg. comp. time | #10 | #100 |
|-------------|-----------------|-----|------|
| DTM KW 1k   | 0.177 sec.      | 24  | 5    |
| DTM KW 5k   | 0.355 sec.      | 116 | 5    |
| DTM KW 10k  | 0.511 sec.      | 28  | 7    |
| DTM KW 12k  | 0.594 sec.      | 27  | 7    |
| DB KW       | 0.267 sec.      | 23  | 0    |
| DTM SUB 1k  | 0.195 sec.      | 13  | 1    |
| DTM SUB 5k  | 0.352 sec.      | 65  | 0    |
| DTM SUB 10k | 0.484 sec.      | 18  | 1    |
| DTM SUB 12k | 0.534 sec.      | 24  | 2    |
| DB SUB      | 0.413 sec.      | 706 | 5    |

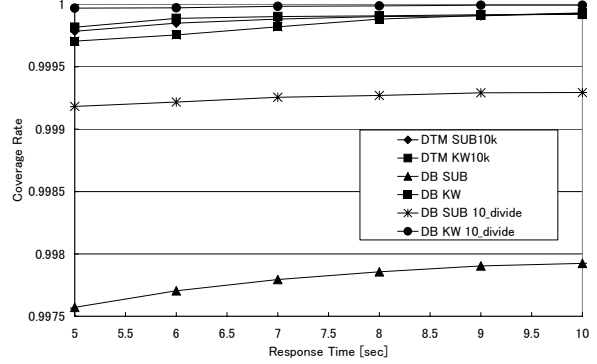


Figure 14: Computation Times for KEYWORD Divided into 10 Tables

The reason why DB requires so much computation time for certain categories is the self-join of the table KEYWORD  $V$  which contains 193,185,919 records. Therefore, we divided the table into multiple tables *in preprocess*. Let  $id(V)$  be a function which returns a set of document IDs in a table  $V$ . We divide  $V$  into multiple tables  $V_i$  which satisfies  $id(V) = \bigcup_i id(V_i)$  and  $id(V_i) \cap id(V_j) = \emptyset$  for  $i \neq j$ . Since the SQL query (3) contains  $V^{(a)}.id = V^{(b)}.id$  in its WHERE clause, the following SQL query can avoid joining tables that do not contain the same documents IDs.

$$\alpha[c, g^{(1)}, 'count'](\sigma_{T:'cancer'}(D)) = value \bowtie value, sum(count) \left( \bigcup_i R_i \right), \quad (4)$$

where  $R_i = value \bowtie value, count(distinct id) count(V_i')$ ,  $\bigcup$  represents UNION ALL operation, and each  $V_i'$  is calculated by the SQL query (3).

Figure 14 shows the results when we compared the aggregation with KEYWORD divided into 10 tables with DTM and the aggregation using a single table for KEYWORD. By dividing the table, the coverage rate within 10 seconds rises from 99.993% to 99.999% for KW and from 99.79% to 99.93% for SUB. Although these experiments were run with a single computer, we can easily run on multiple computers, because such commercial database systems support parallelization.

Figure 15 shows the average computation time for the number of documents for all of the categories. The computation time is observed to be proportional to the number of documents. The high scalability of our method has been confirmed for the amount of data.



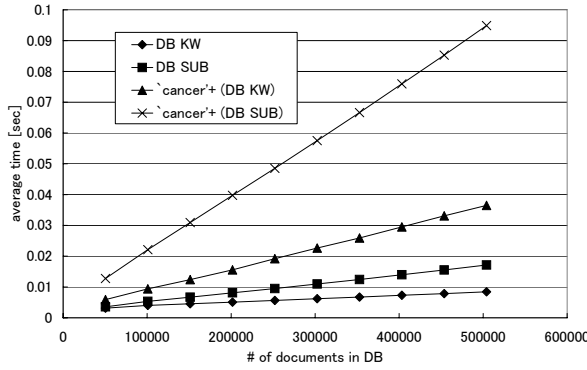


Figure 15: Average Computation Time for the Number of Documents

## 6. DISCUSSION

A top N ranking query often returns a trivial result, which contains only frequent terms in the documents. To measure more informative terms that could show strong relevance to a given subset of documents, the following relative frequency is used [27]. This measure compares the current document subset to the initial document set. Assume  $D$  is the initial document set. A selection operation due to query  $c : v$  returns  $D_s$ . The relative frequency for a term  $p_{jk} = (c_j : v_{jk})$  in the document set  $D_s$  is calculated as

$$relative\_frequency(p_{jk}, D_s) = \left( \frac{c(p_{jk}, D_s)}{|D_s|} \right) / \left( \frac{c(p_{jk}, D)}{|D|} \right),$$

where  $c(p_{jk}, D)$  is the number of documents that contain the term  $p_{jk}$  in the set  $D$ . For example in Figure 11 in Section 5.2,  $relative\_frequency((c_3 : v_{31}), D_s)$  is  $\frac{2}{3} / \frac{4}{10} = \frac{5}{3}$ , where  $D_s$  is  $\{d_2, d_6, d_{10}\}$ .

It is very critical to achieve query response times that are as fast as possible for interactively analyzing a huge amount of text data. A key strategy for speeding up to aggregate the data is to use indexing technology. As mentioned in Section 4, we use the preorder and postorder to check ancestor-descendant containment in a category hierarchy. If we do not use preorder and postorder in a traversed tree, we need to join the table **CATEGORY**  $n$  times to check where a node  $A$  is an ancestor of a node  $B$ , where  $n$  is the length of a path from the  $A$  to the  $B$ . However, we do not need any join operations to ancestor-descendant containment in the hierarchy. The method to index the tree was proposed in 1982 [7], and it recently draws attention as the method to index XML (eXtensible Markup Language) database and to map XML data into the relational database [11], since each XML document is modeled as a DOM (Document Object Model) tree. Several methods such as prefix label [6], Dewey order [26], prime label [28], VLEI code [14], embedding into a k-ary tree [15] are used to index XML. A disadvantage of the methods such as preorder-postorder method and prime label is to need a re-assignment of preorder and postorder of nodes when inserting some nodes into a tree. Since each node has the same label as a prefix of its children in the methods such as prefix label and Dewey order, they do not need to be reassigned the labels when inserting some nodes. However, because they need to compute functions to process string to check ancestor-descendant containment, they need more computation time than the preorder-postorder method. Since we assume that a category hierarchy is far

less frequently updated than new records are inserted into the table  $V$ , we used the preorder-postorder method to index the category hierarchy.

Our data representation is similar to the bag-of-words approach [16], although each term is assigned categories. In the bag-of-words approach, the following sentences are treated as the same content: “(a) X did fail”, “(b) X did not fail”, and “(c) Did X fail?” [19]. Besides the negation and interrogative mood, some auxiliary verbs such as “can” and some verbs such as “want” often indicate the author’s communicative intentions. It is important to associate communicative intentions with predicates by analyzing grammatical features and lexical information. “fail” in the previous examples (a) to (c) are assigned categories and stored in  $R$  as (a) complaint:fail, (b) commendation:not\_fail, and (c) question:fail. These distinctions are instrumental in facilitating problem detection and workload reduction for analysts at customer help centers, for example.

Some papers such as [17] and [18] proposed OLAP systems to analyze a set of documents. The following MDX (Multi-Dimensional eXpression) query which was used in [17] finds all the documents which contain a term “forests” and are published in New York in the first quarter of 1998.

```
SELECT not empty [DocId].members on rows,
      {[Measure].[Tf]} on columns
FROM docInfo
WHERE ([Term].[forest], [1998][quarter 1],
      [location].[New York])
```

In the query,  $[Term].[forest]$  means that a depth of a hierarchy in a TERM dimension is 2, and term “forests” is a child node of “TERM” corresponding to  $\top$  in the TERM dimension. Hierarchies that the existing OLAP systems for texts assume are so simple that it is difficult to integrate the hierarchies with a complex ontology with a huge set of nodes. Fagin et al. proposed Multi-Structural Database (MSDB) to support efficient analysis of large, complex data set. Our method provides a user interactive analysis, although MSDB can semi-automatically segment the data by using analytics operators among some correlated dimensions. In addition, we demonstrated our method by much larger set of data than those of MSDB.

Other papers proposed an OLAP system to analyze a set of documents [25, 24]. Since sales figures mentioned in multiple documents (newspapers) would be a fact in the papers, operations in the system is similar to the conventional OLAP system for structured data. An example operation in the system is to analyze the average sales per product and year. In our case, since a document would be a fact, we can find a change in the number of documents with time. For example, in a call center in a company, call takers make reports of each call by typing in customer information such as name and phone number, selecting call categories such as “technical QA” and typing in brief descriptions of questions or messages from the customer and brief descriptions of answers and/or actions taken. The brief descriptions are written in natural language. The manager of the call center wants to improve productivity, reduce cost, improve customer satisfaction, etc. For example, in a large number of documents related to customers’ calls, he or she would like to find what kinds of topics have recently been increasingly mentioned and which product is associated with specific topics, so that we can take appropriate actions for the improvement of call center productivity and product quality, or create a FAQ (frequently asked question) database.

Our method relies on natural language processing and information extraction, and to build a general purpose extraction module is not a trivial task. IBM makes UIMA (Unstructured Information Management Architecture, <http://uima-framework.sourceforge.net/>) available as an open source software development kit [9]. The UIMA framework is an open, scalable and extensible platform for building analytic applications or search solutions that process text or other unstructured information. It enables developers to build analytic modules, *e.g.* entity extraction modules. It can also enable us to extract various technical and meaningful terms to analyze a huge set of documents by composing modules from multiple annotator providers.

## 7. CONCLUSION

In this paper, we proposed a data representation and its algebra operations to integrate ontologies with OLAP systems to analyze a huge set of textual documents. By using our method, two types of information (structured and unstructured information) can mutually enhance information discovery and analysis capability. The proposed method was implemented with a persistent store using preorder and postorder in a hierarchy. The efficiency of our approach has been confirmed with respect to the computation time. Our method is so efficient and robust that it enables an analyst to interactively analyze a large amount of text data for more context-oriented analysis and decision-making.

## 8. ACKNOWLEDGMENTS

We thank members of text mining group at Tokyo Research Laboratory, IBM Japan for their supports and advices.

## 9. REFERENCES

- [1] S. Agarwal et al. On the Computation of Multi-dimensional Aggregates. *Proc. of Int'l Conf. on Very Large Data Bases*, pp. 506–521, 1996.
- [2] P. Baumann et al. Spatio-Temporal Retrieval with RasDaMan. *Proc. of Int'l Conf. on Very Large Data Bases*, pp. 746–749, 1999.
- [3] P. Bonatti et al. An Ontology-extended Relational Algebra. *Proc. of Int'l Conf. on Information Reuse and Integration*, pp. 192–199, 2003.
- [4] V. Chakaravarthy et al. Efficiently Linking Text Documents with Relevant Structured Information. *Proc. of Int'l Conf. on Very Large Data Bases*, pp. 667–678, 2006.
- [5] E. Charniak. *Statistical Language Learning*. The MIT Press, Cambridge, Massachusetts, 1996.
- [6] E. Cohen et al. Labeling Dynamic XML Trees. *Proc. of SIGACT-SIGMOD-SIGART Symp. on Principles of Database Systems*, pp. 271–281, 2002.
- [7] P. Dietz. Maintaining Order in a Linked List. *Proc. of Symp. on Theory of Computing*, pp. 122–127, 1982.
- [8] R. Fagin et al. Multi-Structural Databases. *Proc. of ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems*, pp. 13–15, 2005.
- [9] D. Ferrucci & A. Lally. Building an Example Application with the Unstructured Information Management Architecture. *IBM Systems Journal*, 43(3):455–475, 2004.
- [10] S. Goil & A. N. Choudhary. High Performance Multi-dimensional Analysis of Large Datasets. *Proc. of Int'l Wks. on Data Warehousing & OLAP*, pp. 34–39, 1998.
- [11] T. Grust. Accelerating XPath Location Steps. *Proc. of SIGMOD Int'l Conf. on Management of Data*, pp. 109–120, 2002.
- [12] H. Gupta et al. Index Selection for OLAP. *Proc. of Int'l Conf. on Data Engineering*, pp. 208–219, 1997.
- [13] M. Gyssens & L. Lakshmanan. A Foundation for Multi-dimensional Databases. *Proc. of Int'l Conf. on Very Large Data Bases*, pp. 106–115, 1997.
- [14] K. Kobayashi et al. VLEI code: An Efficient Labeling Method for Handling XML Documents in RDB. *Proc. of Int'l Conf. on Data Engineering*, pp. 386–387, 2005.
- [15] Y. Lee et al. Index Structures for Structured Documents. *Proc. of Int'l Conf. on Digital Libraries*, pp. 91–99, 1996.
- [16] C. Manning & H. Schütze. *Foundations of Statistical Natural Language Processing*. The MIT Press, 1999.
- [17] M. McCabe et al. On the Design and Evaluation of a Multidimensional Approach to Information Retrieval. *Proc. of Int'l SIGIR Conf. on Research and Development in Information Retrieval*, pp. 363–365, 2000.
- [18] J. Mothe et al. DocCube: Multi-Dimensional Visualization and Exploration of Large Document Sets. *Journal of the American Society for Information Science and Technology*, 54(7):650–659, 2003.
- [19] T. Nasukawa & T. Nagano. Text Analysis and Knowledge Mining System. *IBM Systems Journal*, 40(4):967–984, 2001.
- [20] J. Nielsen. *Designing Web Usability: The Practice of Simplicity* New Riders Publishing, 2000.
- [21] T. Niemi et al. Logical Multidimensional Database Design for Ragged and Unbalanced Aggregation. *Proc. of Int'l Wks. on Design and Management of Data Warehouses*, pp. 7, 2001.
- [22] T. Pedersen & C. Jensen. Multidimensional Data Modeling for Complex Data. *Proc. of Int'l Conf. on Data Engineering*, pp. 336–345, 1999.
- [23] T. Pedersen & C. Jensen. Multidimensional Database Technology. *IEEE Computer*, 34(12):40–46, 2001.
- [24] J. Pérez et al. A Relevance-extended Multidimensional Model for a Data Warehouse Contextualized with Documents. *Proc. of Int'l Wks. on Data Warehousing and OLAP*, pp. 19–28, 2005.
- [25] J. Pérez et al. IR and OLAP in XML Document Warehouses. *Proc. of European Conf. on IR Research*, pp. 536–539, 2005.
- [26] I. Tatarinov et al. Storing and Querying Ordered XML using a Relational Database System. *Proc. of SIGMOD Int'l Conf. on Management of Data*, pp. 204–215, 2002.
- [27] N. Uramoto et al. A Text-Mining System for Knowledge Discovery from Biomedical documents. *IBM Systems Journal*, 43(3):516–533, 2004.
- [28] X. Wu et al. A Prime Number Labeling Scheme for Dynamic Ordered XML Trees. *Proc. of Int'l Conf. on Data Engineering*, pp. 66–78, 2004.