

# A New Approach to Compute Structural Statistics using Keywords in Databases

Ambily Jose<sup>1</sup>, Jobin George<sup>2</sup>, Liliya T Jose<sup>3</sup>

<sup>1</sup>Dept of Computer Science and Engineering,  
BTL Institute of Technology, Bangalore, India.

<sup>2</sup> Software Engineer, Bangalore, India.

<sup>3</sup>Dept of Computer Science and Engineering,  
Vimal Jyothy Engineering College, Kerala, India.

**Abstract**— Query using keywords is one of the user friendly and widely used forms of querying in recent days. The existing systems focused on fetching all or a fixed number of top tuples, which may directly contain the portion of keywords or interconnected to other related tuples via foreign key. These systems will not give any valuable information more than individual interconnected tuple structures.. In this paper we focus on compute structural statistics for keyword by extending the group & aggregate framework. In RDBMS, rows with same values on specific attributes can be grouped together and aggregate functions can be used to aggregate the values of the same group. Tuples are considered as rooted subgraphs, which represents an interconnected tuple structure among tuples. The keywords are separated as dimensional keywords and general keywords. The dimensions of the rooted subgraphs are determined by dimensional keywords and the general keywords are used to compute scores. Based on the score computed for every group, an aggregate function is used to compute the aggregates.

**Keywords** - Keyword search, relational database, Structural statistics.

## I. INTRODUCTION

Search based on keyword in database systems is one of the commonly used research topic. It helps the users to get results which match the keywords. Most of the available studies are based on finding the related values based on the values on the table and values on other tables which is referred by a foreign key reference[1]-[3]. These systems may give a large set of possible results for the given keywords, and sometimes users may struggle to filter the required information. In this paper, we are going to evaluate how to compute statistics on the interconnected tuple- structures using keywords instead of finding interconnected structures among tuples.

## II. STRUCTURAL STATISTICS

Let  $GS = \{R_1, R_2, \dots\}$  be a relational database schema. Here,  $R_i$  in  $GS$  is a relation schema with a set of attributes. Keyword search is allowed on any text-attributes. A relation schema may have a primary key and there are foreign key references defined in  $GS$ . We use  $R_i \rightarrow R_j$  to denote that there is a foreign key defined on  $R_i$  referring to the primary key defined on  $R_j$ . A relation on relation schema  $R_i$  is an instance of the relation schema (a set of tuples) conforming to

the relation schema, denoted  $r(R_i)$ . A relational database (RDB) is a collection of relations.

We model an RDB over  $GS$  as a directed database graph  $GD(V, E)$ . Here we use two types of labeled nodes tuple nodes  $V_t$  and attribute nodes  $V_a$ ,  $V = V_t \cup V_a$  for  $V_t \cap V_a = \Phi$ . A tuple-node in  $r(R')$  labeled with its relation name  $R'$  and  $t_i$  represents a tuple in  $r(R')$ . An attribute-node is labeled  $A_j : a_k$  where  $A_j$  is a text-attribute name and  $a_k$  is an attribute value. Consequently,  $E = E_{tt} \cup E_{ta}$  consists of two types of edges.  $E_{tt}$  is the set of edges among tuple-nodes in  $GD$ , and an edge  $t_i \rightarrow t_j$  in  $E_{tt}$  indicates that there exists a foreign key reference from  $t_i$  to  $t_j$  in RDB.  $E_{ta}$  is the set of edges from tuple-nodes to attribute-nodes. There exists an edge from a tuple-node  $t_i$  to an attribute-node  $A_j : a_k$  in  $E_{ta}$ , if the tuple  $t_i$  has the attribute value  $a_k$  in the attribute  $A_j$  in RDB. Here  $V(G)$  and  $E(G)$  denote the set of nodes and the set of edges of  $G$ , respectively.

If  $k_i$  is contained in either the attribute name or the attribute value in the node label then we can say that there exist an attribute-node  $v$  contains a keyword  $k_i$ . We also say a tuple-node  $u$  contains a keyword  $k_i$  if there is a path from tuple node  $u$  to an attribute node  $v$  contains  $k_i$ .

**Virtual Tuple.** A virtual tuple is a tree representation of the maximum subgraph at a tuple-node  $t_\gamma$  in  $GD$  and denoted  $e$  as  $V_{tuple}$ , or explicitly  $V_{tuple}(t_\gamma)$  if it is rooted at a tuple-node  $t_\gamma$ . The maximum subgraph of  $t_\gamma$  is the induced subgraph of all nodes that are reachable from  $t_\gamma$  in  $GD$ . The tree representation of  $V_{tuple}(t_\gamma)$  is, if a node  $u$  links to a node  $w$  which is predecessor of  $u$ , the edge  $(u, w)$  will be deleted. Otherwise, if two nodes  $u$  and  $v$  link to a node  $w$  in  $GD$ , we create an additional copy  $w'$  of  $w$ , and let  $u$  link to  $w$  and  $v$  link to  $w'$ . All leaf nodes in  $V_{tuple}(t_\gamma)$  must be attribute-nodes. A  $V_{tuple}(t_\gamma)$  includes all information a tuple-node  $t_\gamma$  can reach.

**Dtree, Gtree, and DGtree.** Given a set of keywords  $\{k_1, k_2, \dots\}$ , a  $Dtree(t_\gamma)$  is a minimal subtree of  $V_{tuple}(t_\gamma)$  which contains all dimensional-keywords by connecting to the attribute-nodes that contain the given dimensional keyword. A  $Gtree(t_\gamma)$  is also a subtree of  $V_{tuple}(t_\gamma)$  by removing all the subtrees rooted at a tuple node that do not have any attribute-node containing general keywords. A  $Gtree(t_\gamma)$  matches a query if it contains at least one general-keyword. A  $Gtree(t_\gamma)$  may not contain all the general-keywords. Given a set of

keywords, there exists one  $Gtree(t\gamma)$ , but many distinctive  $Dtree(t\gamma)$ , by definition. A  $DGtree(t\gamma)$  consists of two parts, where the first part is  $Dtree(t\gamma)$  and the second part is  $Gtree(t\gamma)$ .

In this paper, we study a new structural statistics keyword query  $Q = (Qd, Qg, \alpha, \beta)$  against an RDB over GS. It consists of two sets of keywords, namely  $Qd$  and  $Qg$ , for  $Qd \cap Qg = \Phi$ , a score function  $\alpha$ , and an aggregate function  $\beta$ . We call a keyword in  $Qd$  and  $Qg$  as a dimensional-keyword and a general-keyword, respectively. The two sets of keywords,  $Qd$  and  $Qg$ , together specify a set of trees  $T$  to be computed. A  $DGtree$   $T_i$  in  $T$  with root  $t\gamma = root(T_i)$  consists of two subtrees, a  $Dtree$  rooted at  $t\gamma$  for  $Qd$ , and a  $Gtree$  rooted at the same  $t\gamma$  for  $Qg$ , denoted as  $Dtree(T_i)$  and  $Gtree(T_i)$ , respectively. The set of trees,  $T$ , are grouped into different groups. Let  $T_i$  and  $T_j$  be two trees in  $T$ .  $T_i$  and  $T_j$  belong to the same group if  $Dtree(T_i)$  is isomorphic to  $Dtree(T_j)$ . Here score function  $\alpha$  to be any possible algebraic function based on TF-IDF, namely,  $tf_w(T)$  and  $idf_w$ . Let  $T = \{T_1, T_2, \dots\}$  be a set of  $DGtrees$  in the same group. We consider every  $Gtree(T_i)$  for  $T_i \in T$  as a virtual document, by merging all attribute names and attribute values in the tree into a multiset. Then,  $tf_w(T_i)$  is the number of times the keyword  $w \in Qg$  appears in the corresponding virtual document, and  $idf_w$  is calculated as follows[5]:

$$idf_w = \frac{|T|}{df_w(T)} \quad (1)$$

where  $|T|$  is the number of  $Gtrees$  in  $T$  in the group, and  $df_w(T)$  is the number of  $Gtrees$  that contain the keyword  $w \in Qg$  in the group. The tree level ranking function[4] is such an algebraic function based on  $tf_w(T)$  and  $idf_w$ . The  $\alpha$  function is to be applied to  $Gtree(T_i)$  for  $T_i \in T$  to give such  $Gtree(T_i)$  a score. Factors that can be involved in an  $\alpha$  function can be, for example, to give a high score for a term if it is close to the root. For efficiency consideration we require that factors in an  $\alpha$  function for a tree must be computable from the factors of its subtrees. The aggregate function  $\beta$  aggregates the scores computed for  $DGtrees$  in the same group. An aggregate function can be any SQL aggregate functions (min, max, sum, avg, and count). The output for the group  $T$  is  $(TA, \omega)$  where  $TA$  represents the  $Dtree$  for the group, and  $\omega = \beta(\{\alpha(Gtree(T_1)), \alpha(Gtree(T_2)), \dots\})$ .

### III. SOLUTION OVERVIEW

Given a structural statistics keyword query  $Q = (Qd, Qg, \alpha, \beta)$  over an RDB, a naive solution is impractical for the following two main reasons. 1) The number of possible  $Dtrees$  and  $Gtrees$  can be very large. It is infeasible to compute. It is worth noting that all the existing solutions focus on finding top-k answers if they make use of ranking and allow some (not all) keywords to be contained in the answers. 2) It is costly to group  $Dtrees$  into groups even though tree isomorphism checking is polynomial.

In this paper, to compute a structural statistics keyword query,  $Q$ , we propose a two-step approach (Algorithm 1). In the first step, we generate a set of label-trees for dimensional keywords (LDs), denoted as  $L = \{LD_1, LD_2, \dots\}$ , such that every  $DGtree$  to be computed will conform to a unique LD.

In the second step, we compute the structural statistics keyword query  $Q$  using  $L$ . In brief, for every  $LD_i$ , we compute all  $DGtrees$ , denoted as  $T_i = \{T_{i1}; T_{i2}; \dots\}$  that conform to  $LD_i$ , group all the trees in  $T_i$  into groups based on  $Dtree(T_{ij})$  and compute  $\alpha(Gtree(T_{ij}))$  for every  $T_{ij} \in T_i$ , and then compute  $\beta$  for every group. The main idea behind label-trees is to avoid tree isomorphism checking and enumerating all possible  $DGtrees$ .

The algorithms to generate all label-trees and to compute structural statistics keyword query using the label-trees will be discussed in the following sections.

#### ALGORITHM 1. Structural-Statistics(GS, Q, GD)

```
1: L = LD-Gen(GS, Q);
2: for every  $LD_i \in L$  do
3: LD-Eval(Q,  $LD_i, G_D$ );
```

First, we define a label-graph,  $G_s(V, E)$ , for keyword search for a database schema  $GS$ . Here,  $V$  is a set of nodes such that  $V = V_R \cup V_A$ , where  $V_R$  is a set of nodes labeled with relation names, called relation nodes, and  $V_A$  is a set of nodes labeled with attribute names for those text-attributes only.  $E$  is a set of edges such that  $E = E_{RR} \cup E_{RA}$ . An edge  $R_i \rightarrow R_j$  appears in  $RR$  if there is a foreign key reference from  $R_i$  to  $R_j$ . An edge  $R_i \rightarrow A_j$  appears in  $E_{RA}$  if  $R_i$  has an attribute  $A_j$ . Second, we define a rooted tree for every relation  $R_i$  in  $G_s$ , denoted as  $LV(R_i)$ , which is a labeled tree for all virtual tuples rooted at  $t_\gamma \in r(R_i)$ .  $LV(R_i)$  is a connected tree representation of the maximum subgraph rooted at  $R_i$  in  $G_s$ . The tree representation is done as follow. In  $LV(R_i)$ , if a node  $u$  links to a node  $w$  which is an ancestor of  $u$ , the edge  $(u, w)$  will be deleted. Otherwise, if two nodes  $u$  and  $v$  link to a node  $w$  in  $GS$ , we create an additional copy  $w'$  of  $w$ , and let  $u$  link to  $w$  and  $v$  link to  $w'$ .

Consider a structural statistics keyword query  $Q = (Qd, Qg, \alpha, \beta)$  against the data graph  $G_D$  using the label-graph  $G_s$ . We further give a specific  $LV$  for computing  $Dtrees$ . A label-tree for dimensional-keywords  $LD_i(R_j)$  is a subtree of  $LV(R_j)$  rooted at  $R_j$  that contains at most  $|Qd_j|$  attribute-nodes as leaf nodes. The attribute-nodes in  $LD_i(R_j)$  possibly contain all the dimensional-keywords in  $Qd$ .

#### IV. GENERATE ALL LDs

The solution is as follows, first, we precompute all  $LVs$  for  $G_s$  because the set of  $LVs$  is query independent. The algorithm to compute all  $LVs$  is shown in Algorithm 2 [5]. For each node  $R$  in  $GS$ , we calculate  $LV(R)$  using a breadth first search from  $R$  in  $GS$  until all nodes that can reach from  $R$  are added into  $LV(R)$ . For node  $u$  that is visited more than once in the breadth-first search, we create an extra copy in  $LV(R)$  if  $u$  is not visited from its descendant. Second, in order to efficiently generate all LDs for all possible dimensional-keywords, we construct an inverted index, called the dimensional inverted index (DII), using the names and values of the attributes in the RDB. The inverted index helps to find the attributes in a relation that a dimensional-keyword  $di$

matches. In detail, for each possible dimensional-keyword,  $w$ , in DII, there is a list of entries to describe the attributes in a relation the keyword  $w$  matches. We denote the list for  $w$  as  $\text{list}(w)$ . Each entry  $e \in \text{list}(w)$  has three fields:  $e = \{\text{Type}, \text{Rel}, \text{Attr}\}$ . Here, Type can be either Name or Value. When Type = Name, it means  $w$  is an attribute name. When Type = Value, it means  $w$  is an attribute value. Rel and Attr indicate the relation and the attribute that  $w$  matches.

#### ALGORITHM 2. LV-Gen(GS)

Input: The label-graph GS (VR U VA; ERR U ERA )

Output: The set of all LVs.

```

1: S ← Φ;
2: for all R ∈ VR do
3: add R with links to all its text-attributes into LV(R);
4: Q ← Φ;
5: Q:enqueue(R);
6: while Q ≠ Φ do
7:   R1 ←Q.dequeue();
8:   for all R1 → R2 ∈ E(Gs) do
9:     let R'₂ be a copy of R₂;
10:    add an edge in LV(R) from R1 to R'₂
11:    add links from R'₂ to all its text-attributes;
12:    Q:enqueue(R'₂);
13: S ←S U {LV(R)};
14: return S

```

We design a new algorithm to generate all LDs using  $\text{list}(d_i)$  for  $d_i \in Q_d$ , as shown in Algorithm 3. The algorithm[5] generates all LDs for a structural statistics keyword query. First, it collects information if a keyword  $d_i$  matches the relation nodes in each LV (lines 2-6). Given the set of LVs  $S$ , for each LV in  $S$ , we use  $\text{LV.list}_i$  to maintain the set of candidate entries in  $\text{list}(d_i)$  that may contribute to generating LDs from LV. Second, in a for loop (lines 7-10), for each combination of nodes  $e_1, \dots, e_{|Q_d|}$  that contain keywords  $d_1, \dots, d_{|Q_d|}$ , respectively, in a certain  $\text{LV} \in S$ , we generate a LD by constructing a minimum tree that contains nodes  $e_1 \dots e_{|Q_d|}$  in LV. Given LV and  $e_1, \dots, e_{|Q_d|}$ , the LD is unique and can be computed as follows: for each leaf node of LV, we remove all the leaf nodes that do not belong to  $\{e_1 \dots e_{|Q_d|}\}$ . We do this iteratively until no leaf node is removed. The result is a minimum tree that contains nodes  $e_1, \dots, e_{|Q_d|}$  in LV

#### ALGORITHM 3. LD-Gen(GS, Q, S, DII)

Input: The label-graph GS, a query  $Q = (Q_d, Q_g, \alpha, \beta)$ , the dimensional inverted index DII and the set of LVs,  $S$ .

Output: The set of all LDs.

```

1: Q ← Φ;
2: for each keyword  $d_i$  in  $Q_d$  do
3:   for each  $e \in \text{list}(d_i)$  do
4:     for each  $\text{LV} \in S$  do
5:       if  $e.\text{Rel} \in \text{LV}$  then
6:          $\text{LV.list}_i \leftarrow \text{LV.list}_i \cup \{e\}$ ;

```

7: for each  $\text{LV} \in S$  do.

8: if  $\text{LV.list}_i \neq \Phi$  for any  $1 \leq i \leq |Q_d|$ . then

9: for all  $(e_1, \dots, e_{|Q_d|}) \in \text{LV.list}_1 \times \dots \times \text{LV.list}_{|Q_d|}$  do

10: construct a minimal LD that only contains attribute-nodes  $(e_1, \dots, e_{|Q_d|})$ ;

11:  $Q \leftarrow Q \cup \{\text{LD}\}$ ;

12: return  $Q$ ;

#### V. EVALUATE ALL LDS

Here, for a structural statistics keyword query  $Q = (Q_d, Q_g, \alpha, \beta)$ , we give a two-phase approach to compute structural statistics for all the groups under a given LD.

##### A. Two-Phase Approach

We propose a new two-phase approach, namely, bottom-up followed by top-down, after pruning unnecessary nodes/ edges from GD that need to evaluate an LD(R). For a given LD(R) to be evaluated, let LV(R) be the labelled LV that generates LD(R). In the bottom-up phase, we collect statistics for TF-IDF, namely,  $\text{tf}_w(T)$  and  $\text{idf}_w$  using general-keywords. The statistics collection is done in a bottom-up fashion, and will finish when it finally evaluates the relation R which is the root of LD(R). Let  $R'$  be the set of tuples in R that contain statistics for at least one general-keyword in  $Q_g$ . At the end of this phase, we have a set of  $\{\text{Gtree}(t_\gamma)\}$  for  $t_\gamma$  in  $R'$ , and we compute all  $\alpha(\text{Gtree}(t_\gamma))$ . In the top-down phase, we start from those tuples  $t_\gamma$  in  $R'$ . We retrieve all leaf attribute nodes that  $t_\gamma$  can reach using the depth first search from  $t_\gamma$ . If the leaf attribute nodes contain all the required dimensional keywords in  $Q_d$ , then  $t_\gamma$  has a valid Dtree( $t_\gamma$ ). At the end of this top-down process, we compute  $\gamma$  for all such  $t_\gamma$  in  $R'$  that has a valid Dtree( $t_\gamma$ ).

We first discuss the pruning process using dimensional keywords before the two-phase approach. Given a dimensional keyword  $d_i$ , a node  $v$  in  $G_D$  can be pruned by  $d_i$  if there does not exist a Dtree that contains both keyword  $d_i$  and node  $v$ . We call  $v$  an unnecessary node if  $v$  can be pruned by any dimensional keyword in the query. Given an LD, we consider whether a keyword  $d_i \in Q_d$  has enough pruning power using  $\text{Pow}(d_i, \text{LD})$  which is computed as follows[5]: suppose  $X(A, B)$  is the join selectivity in relation A that can join a tuple in relation B, i.e., the fraction of tuples of A that have a matching tuple in B. Suppose the attribute in relation  $R_i$  in LD contains  $d_i$ , and assume  $P_i$  is the path from the root of LD to  $R_i$ . We have

$$\text{Pow}(d_i) = \frac{1}{R_i.c(d_i) \cdot \prod_{(A,B) \in P_i} X(A,B)} \quad (2)$$

Here  $R_i.c(d_i)$  is the number of tuples in  $R_i$  that contain the keyword  $d_i$  in the specified attribute of  $R_i$  in LD. The basic idea behind is as follows: for any join sequence  $A \bowtie B \bowtie C$ , we assume that for any tuples  $a \in A$ ,  $b \in B$ , and  $c \in C$ , whether  $a$  can be joined with  $b$  and whether  $b$  can be joined with  $c$  are independent to each other.  $R_i.c(d_i) \cdot \prod_{(A,B) \in P_i} X(A,B)$  is the expected number of root nodes to reach a node with keyword  $d_i$ . Intuitively, a keyword that ends up with a smaller expected number of the root nodes has higher pruning power.

For any general-keyword  $g_j \in Q_g$ , suppose for any text-attribute  $A_i$ ,  $R(A_i)$  is the relation of  $A_i$ , and  $P(A_i)$  is the path from the root of LD to the relation  $R(A_i)$ . We compute the pruning power of a general-keyword as follows[5]:

$$\text{Pow}(g_i, LD) = \frac{1}{\sum_{A_i \in LD} R(A_i).C(g_i) \cdot \prod_{(A,B) \in P(A_i)} X(A,B)} \quad (3)$$

Based on the equations, we decide whether a dimensional keyword  $d_i \in Q_d$  has enough power to prune, or in other words, it is cost-effective to reduce  $G_D$ . We sort all dimensional-keywords  $d_i \in Q_d$  in decreasing order. If the pruning power of  $d_i$  is larger than the largest pruning power of all general-keywords in  $Q_g$ , then we use  $d_i$  to reduce  $G_D$  by removing all the tuple-nodes that cannot reach any attribute nodes containing  $d_i$ .

The algorithm is shown in Algorithm 4[5]. We assume that the structure of the data graph is held in memory. First, we prune unnecessary nodes from  $G_D$  using  $Q_d$  if they have enough pruning power (lines 5-7). Second, in the bottom-up phase, we compute trees in a sense to collect all needed information to compute  $\alpha$  for every Gtree using  $Q_g$ . It is done from the leaf toward the root which is a tuple in  $r(R)$  for the  $LV(R)$  that generates the given LD (lines 9-25). Finally, in a top-down phase, we aggregate for each group based on  $Q_d$  (lines 27-31).

#### ALGORITHM 4. LD-Eval(Q, LD, GD)

Input: A query  $Q = \{Q_d, Q_g, \alpha, \beta\}$ , LD and a database graph  $GD(V, E)$ .

Output: Aggregates for all groups

```

1:  $\Gamma \leftarrow \Phi$ ;
2: let  $LV(R)$  be the LV that generates LD;
3: for all relation node  $P \in LV(R)$  do
4:  $P.set \leftarrow \Phi$ ;
5: for all  $d_i \in Q_d$  sorted by decreasing order
   of  $\text{Pow}(d_i, LD)$  do
6:   if  $\text{Pow}(d_i, LD) > \max_{g_i \in Q_g} (\text{Pow}(g_i, LD))$  then
    $GD \leftarrow \text{prune}(GD, d_i, LD)$ ;
8: // The bottom-up phase
9: for all relation node  $P \in LV(R)$  in the order from
   leaves to the root do
10:  for  $i=1$  to  $|Q_g|$  do
11:   for all tuple-node  $t_p \in P$  contain( $g_i$ ) do
12:     $t_p.cnt_i \leftarrow t_p.cnt_i + t_p.count(g_i)$ 
13:     $P.has_i \leftarrow P.has_i \cup \{t_p\}$ ;
14:     $P.set \leftarrow P.set \cup \{t_p\}$ ;
15: for all child of  $P$ ,  $C$ , in  $LV(R)$  do
16:  for all node  $t_c \in C.set$  do
17:   for all node  $t_p \in P$  such that  $t_p \rightarrow t_c \in E(GD)$  do
18:     $P.set \leftarrow P.set \cup \{t_p\}$ ;
19:  for  $i=1$  to  $|Q_g|$  do
20:    $t_p.cnt_i \leftarrow t_p.cnt_i + t_c.cnt_i$ ;
21:   if  $t_c.cnt_i > 0$  then  $P.has_i \cup \{t_p\}$ 
22: for all tuple-node  $t \in R.set$  do
23:  for  $i=1$  to  $|Q_g|$  do

```

```

24:   $tf_{g_i}(t) \leftarrow t.cnt_i$ ;  $df_{g_i}(LV(R)) \leftarrow |R.has_i|$ ;
25:   $t.score \leftarrow \alpha(t)$  using all  $tf_{g_i}(t)$  and  $df_{g_i}(LV(R))$ ;
26: // The top-down phase
27: for all  $t \in R.set$  do
28:  let  $a_i$  be the attribute value in the attribute
    $A_i \in att(LD)$  that contains  $d_i$ , for all  $d_i \in Q_d$ ;
29:  let  $\gamma$  be a group represented by  $(a_1, a_2, \dots, a_{|Q_d|})$ 
30:   $\gamma.score = \beta(\gamma.Score, t.score)$ ;
31:   $\Gamma \leftarrow \Gamma \cup \{\gamma\}$  if  $\gamma$  not in  $\Gamma$ ;
32: return  $\Gamma$ 

```

## VI. CONCLUSIONS

In this paper, we studied how to compute structural statistics for all groups of tuples in RDB using the keyword query. By using a new two-step approach, we compute all label-trees for a structural statistics keyword query, and a new two-phase algorithm used to compute group-&-aggregate using the label-trees computed.

## REFERENCES

- [1] S. Agrawal, S. Chaudhuri, and G. Das, "DBXplorer: A System for Keyword-Based Search over Relational Databases," Proc. 18th Int'l Conf. Data Eng. (ICDE '02), 2002.
- [2] V. Hristidis, L. Gravano, and Y. Papakonstantinou, "Efficient IRStyle Keyword Search over Relational Databases," Proc. 29th Int'l Conf. Very Large Data Bases (VLDB '03), 2003
- [3] V. Hristidis and Y. Papakonstantinou, "DISCOVER: Keyword Search in Relational Databases," Proc. 28th Int'l Conf. Very Large Data Bases (VLDB '02), 2002.
- [4] Y. Luo, X. Lin, W. Wang, and X. Zhou, "Spark: Top-K Keyword Query in Relational Databases," Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '07), 2007
- [5] Lu Qin, Jeffrey Xu Yu, and Lijun Chang, "Computing Structural Statistics by Keywords in Data bases", IEEE transactions on knowledge and data Engg., oct 2012
- [6] L. Qin, J.X. Yu, and L. Chang, "Computing Structural Statistics by Keyword in Databases," Proc. IEEE 27th Int'l Conf. Data Eng. (ICDE '11), 2011.
- [7] L. Qin, J.X. Yu, and L. Chang, "Keyword Search in Databases: The Power of Rdbms," Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '09), 2009

Ambily Jose received the B.Tech. degree in Computer Science and Engineering from Cochin University of Science and Technology, Cochin, Kerala. At present pursuing the Master of Technology in Computer Science and Engineering Department at BTL institute of Technology, Bangalore.

Jobin George currently working as a Senior Software Engineer in Wipro Technologies, Bangalore, India

Liliya T Jose currently working as an Assistant Professor in Vimal Jyothy Engg. College affiliated to Kannur University, Kerala