

EFFICIENT DATA RETRIEVAL FROM LARGE GRAPH DATABASE FOR MEDICAL APPLICATIONS

ABSTRACT

Graphs can be used to easily represent the relationship between large quantities of data. A graph can also serve as a general-purpose substrate for evaluating decision-making algorithms. Genomic and biological data are complex that can be represented as graphs. Generally, scientific datasets are large, complicated and are unstructured data which require fast processing. Big data refers to the data sets that are very large, unstructured and complex. Data capturing, storing, querying are some of the challenges faced in processing of big data. Traditional data processing applications are difficult to apply on such huge volume of datasets. So, pre-processing and indexing techniques can help to improve the efficiency. Datasets in the medical and scientific fields keep increasing in large amounts due to various research and experiments. In order to provide meaning to this large volume of data, efficient methods to store, retrieve and analyse must be provided so that processing can be made easier. In real applications most of the graph data are noisy and incomplete. So it has become increasingly important to retrieve graphs in the graph database which match the query graph approximately, rather than exact graph matching. Hence, based on a given query graph, similar graphs are retrieved efficiently. It can be used in medical applications like finding the similar disease patterns. The algorithms proposed are evaluated against existing algorithm and are found to be efficient.

Keywords – structure based indexing, graph similarity, frequent fragments, DFS coding

1 INTRODUCTION

A large number of Data in various applications can be represented using graphs. Chemical compounds, biological data, social network etc are some of the examples. Efficient graph data management techniques are required by graph data models. In order to facilitate this, we have built a heuristic graph matching model to efficiently find graph similarity from a graph database i.e. given a database consisting of n graphs where $D=\{g_1, g_2, \dots, g_n\}$ and a query graph q , finding all similar graphs with respect to the query graph. Indexing of graphs are required in order to efficiently search the similar graphs. Path based indexing is one of the existing and widely used method for indexing. It traverses all possible paths and then indexes the graph using hashing. We have proposed a structure based indexing approach which indexes based on the graph structure. In this method, using efficient pruning and retrieving the frequent fragments, the number of possible paths to be traversed can be greatly reduced.

1.1 BIG DATA

The term big data refers to the data sets that are very large, unstructured and complex. The traditional data processing applications are inadequate since the data set are large and complex in big data. Big data includes the data sets whose storage is beyond the typical database software tools ability to capture, analyze, manage and store. Data capturing, data curation, storage, querying and sharing are some of the challenges in big data analytics. Analysing accurately such large sets of data will give efficient insights in the business process modelling and scientific decision making systems which can result in greater operational efficiency, cost reduction and reduced risk.

Three important big data characteristics are volume, veracity and variety. Some of the big data stores are Hbase, Hive, HDFS etc. Map reduce is one of the best known method for turning raw data into useful vital information. Map Reduce is a procedure of taking a large data set and doing computations on the large data set across multiple systems or computers in parallel. Map reduce

serves as a model for how program or a information is often used to refer to the implementation of the model. Map Reduce consists of two different parts. The first is the Map function. Filtering and sorting is done by the Map function. It takes data and places it inside of categories so that it can be analyzed later. The second is the reduce function which combines all the data and provides a summary.

1.2 GRAPH DATABASE

A graph information could be a information which uses linguistics queries using graph structures which has edges, nodes and properties to represent and store information. The system's key idea is that the graph (which can be relationship or edge), that relates directly information things within the store. The relationships permit information within the store to be connected along and is retrieved in one operation in most of the cases. Graph databases use nodes, edges and properties.

1.3 RELATIONAL AND GRAPH DATABASE

The graph databases cannot handle all types of relational queries and there is no natural way to map SQL queries to graph queries. So graph databases cannot considered as a replacement for the relational databases. On the other hand, the relational databases mainly deals with the tables which have right and predetermined schema. So we can write queries up to a particular and fixed depth by joining those normalized tables, but what do we do if the depth of the interrelationship cannot be predetermined for the given application and the input data is arbitrary, changing and ad hoc in nature, in these cases graph databases comes into rescue. It is intuitive to represent the ad hoc data as graph, and we can recursively iterate the formed graph till any depth to analyse, match pattern and for digital learning. In relational database management system the data sets are progressively filtered and grouped, while graphs are usually navigated and recursively defined depth but not pre-determined joins.

Compared to relational database, graph databases are faster in case of associative data sets as the data sets maps more naturally to the graphs as objects. So, the associative data sets that are represented as graphs scale well for large datasets. In most of the object oriented programming scenario it is more intuitive to relate the data sets as graphs. But still we cannot say that graph databases are the general replacement for the relational databases. Graph databases scales well as they don't involve complex joins of different normalized table as in the relational database management system. Graph databases work even faster for graph-like queries, for example given a county map, finding the shortest path between the two cities. So, the graph database has its own advantages and specific applications and it is not considered as the general replacement of the relational data management system.

1.4 GRAPH VISUALISATION

As the scientific data has more interrelated attributes, it is more intuitive to store it in a graph database. In real applications most of the graph data are noisy and incomplete. So it has become increasingly important to retrieve graphs in the graph database which match the query graph approximately, rather than exact graph matching. Simulating the gene function, cataloging biological information, DNA and protein structures mapping are some of the applications. Decision Making is done confidentially by the accuracy of the big data. Greater operational efficiency, reduced risk and cost reductions are meant by the better decisions. The novel work is to propose an efficient data storage and retrieval model which supports effective querying based on storage and to develop an infrastructure to build knowledge based on research applications by a predictive modelling. Disjoint partition based filter, branch based filtering models suffers some disadvantages like finding the optimal disjoint partition. To develop hybrid filtering model, with heuristics for scientific data to effectively find graph similarity search over large graph datasets.

1.5 PROTEIN – PROTEIN INTERACTION(PPI)

The most biological processes in a cell, including gene expression, morphology, cell growth, proliferation, motility, nutrient uptake, intercellular communication and apoptosis are facilitated by the proteins which are the workhorses of the cell. Protein analysis mainly focussed on single/one proteins, but because of the proteins interact with other proteins for majority of functioning, they should be monitored in the context of their interactions partners to fully understand how proteins interact with each other and identify biological networks which is important to understand how protein functions in the cell. PPIs refers to the interactions between proteins as a result of biochemical event of electrostatic forces. Proteins rarely act alone.

The protein in the organism hold special status. Every phenomenon of life goes through these structure and function of protein to be reflexed [8]. The protein interactions are important in studying the inter atomic system of the living cell and helps in the analysis of different diseases like Cancer, Alzheimer's disease, Creutzfeld-Jacob [21]. Protein – protein interactions are studied from different perspectives right from the bioinformatics, quantum chemistry and molecular dynamics. Proteins bind to each other through a combination of van der Waals forces, hydrophobic bonding, and salt bridges at specific binding domains on each protein. Constant development in these fields lead to the large scale growth of protein – protein interaction graphs. Making meaningful information out of this large data sets will empower the current knowledge on disease pathogenesis and biochemical cascades, as well as it provides new therapeutic targets.

1.6 FREQUENT PATTERN MINING

Subsequence, itemsets or substructures are some of the frequent patterns that seem during a knowledge set with frequency number which is less compared to user given value of the threshold. . A substructure will talk over with totally different structural forms, like sublattices, subgraphs, subtrees etc

which can be combined along with subsequences or set of items. If a substructure occurs often during a graph information, it is referred to as a (frequent) structural pattern. Finding frequent patterns plays an important role in associations of mining, correlations, and plenty of alternative fascinating relationships among knowledge. Moreover, it helps in knowledge compartmentalization, classification, clustering, and alternative data processing tasks further. Frequent pattern mining is a very important data processing task and a targeted theme in data processing analysis.

2 LITERATURE SURVEY

2.1 GRAPH DATABASE

Emir Septian Sori Dongoran et.al [2] proposed that not all kinds of information are appropriate for graph database, one that matched is that the molecular graph information kind. It has featured labelled vertices and afloat edges. Graph information even have categorisation method. For molecular data kind study case, [2] proposed a brand new algorithm GraphGrep which is the most acceptable technique as a result of it assume every node within the graph information incorporates a distinctive variety(node-id) and label (node-label). Therefore it is appropriate for molecular datatype. [2] uses a hash table (fingerprint) as an index, examining the graph information fingerprint with graph query fingerprint to filter the database and use Ullman algorithm to perform subgraph matching. Using [2] we have a tendency to additionally get the foremost economical length-path supported the deepest depth in a graph query.

Patil, Shefali et.al [7], it suggested an overall view about the graph database .Data has been hold on in a tabular form thereby increasing the classification and readability. The illustration of information in the sort of graph lends itself well to schema which is dynamic and data in a well structured manner. No customary system or source language has been denied for graph information. [13] Discussed current applications and implementations of graph

databases, giving an summary of the various sorts out there and their application.

L.Zou et.al [9] suggested the conversion of RDF data into large graph database. The RDF is defined as the Resource Description Framework. It is a data model for modeling Web objects as part of developing the semantic web. It has been used in various application (ex-Wikipedia). Querying the large RDF data is very complex. So the RDF database is changed into graph database (RDF graph). The SPARQL query is changed into query graph. Now the similarity search over the graph database is done to provide the results. This problem now has become the sub graph matching over a large graph database. We already have lot of algorithm for the query graph retrieval over large graph database. In this method the RDF data is converted into Adjacency list table. When the query is converted into adjacency list table it will be easy to search over the RDF graph. Many web applications are based on RDF data. The RDF data are getting updated frequently so it is good to use graph database to represent RDF data.

C.Xiao et.al [11] discussed that Graphs are widely used to model complicated data semantics in many applications like social networks, chemistry, pattern recognition, etc. The graph edit distance is suggested in this method. Using the edit distance method the subgraph matching over a large graph database is done. The threshold value is given by the user. The edit distance of the subgraphs less than or equal to the threshold corresponding to the given query graphs will be the output. This method is also the similarity search not the exact match over the graph database. The major drawback from this method is this method follows path based indexing. But this method handles both scattered and clustered edit operations. Using graph database we can accommodate the future increments or updates will be added to the graph database.

DePiero F.W. et.al [23] stated that the graph structure is the very important means which models the schema less and complicated models such as chemical compounds, protein – protein interaction network, road network, chemical compounds and knowledge query systems [14]. Existing prediction algorithms will not work well in these situations. Introducing use of existing information which was collected for information management in high performance computing systems to build accurate coarse-grained prediction systems. Object abstraction can provide the pros of both files and blocks with better abstraction and flexibility.

2.2 GRAPH INDEXING

Xifeng Yan et.al [3] proposed the problems of categorization graphs and propose a unique answer by applying a graph mining technique. Completely different from the prevailing path-based strategies, [16] then proposed a new technique referred to as graph Index that makes use of frequent substructure as the basic categorization feature. Frequent substructures are important candidates since they explore the internal characteristics of the information and are comparatively stable to information updates. To scale back the dimensions of the index structure, 2 techniques, discriminative fragments and size-increasing support constraint are introduced. [22] showed that gIndex has ten times smaller index size, however achieves ten times higher performance as compared with a typical path-based technique.

G. Wang et.al [5], it detailed mainly about structure of the index for the problem of similarity search on a group of enormous distributed graphs and proposes an Q-gram plan which is an economical compartmentalization mechanism. By moldering the graphs into small fragments i.e little grams (which is organized by k-Adjacent pattern of the tree) and pairing-up on those k-Adjacent pattern trees, we can efficiently calculate the lower bound estimation based on their edit distance and can be used for candidate filtering.

Xiaoli Wang et.al [8], it proposed SEGOS, which is one of the query processing framework and method for efficient indexing for graph similarity search. A good two-level index is made off-line supported by graph decomposition subunits initially. Then, a completely unique search strategy supported by the index is projected as in [17]. Two algorithms custom-made from CA and TA methodology are integrated into the projected strategy to boost graph search. Graph pruning is continuously supported by the projected framework which is straightforward to be pipelined to support graph pruning at regular intervals. Based on the two real datasets which are taken as samples, in depth experiments are conducted and the effectiveness and quantifiability of the approaches are tested.

J. Rocha [21] discussed the importance of indexing. Graph is mainly indexed to facilitate the sub graph isomorphism and similarity queries. The project work consisted of various levels of structures [18]. The primary structure is first composed of a directed acyclic graph (DAG) that contains a vertex for each of the unique, induced sub graphs of the graphs of the database. The hash table is the secondary structure, in which cross-indexes each sub graph for fast isomorphic lookup. The key challenge in using inputs of object-oriented applications are no longer expressed as stable file identifiers; rather than they become much more dynamic and unstatic, hidden inside application logic.

F.Bai et.al [24] proposed that querying is easy because of using index structures, but constructing the index structures is costly and resource consuming. The problem boils down for finding all fragments in a large set of graph that matches a given user query graph pattern [42]. Here we can propose a multi-step R-join procedure with fetch step and filter step based on cluster-based join-index with graph codes. The new algorithm outperforms the old traditional algorithm in case of the elapsed time, and it reduces dominant data transmission cost over the network. Since each step involves filter and fetch, the overall complexity of using this method is high. The algorithm is optimized using a stack-based algorithms to handle pattern queries in directed acyclic

graphs, by building index structures for predecessor list which is used for the purpose of easy querying[26]. It is because of the index structures, the querying becomes easy.

2.3 SIMILARITY OF GRAPHS

W.Zheng et.al [1], proposed the study of similarity search of graph, which retrieves all the graphs which are similar to a given query graph under the properties of edit distance of the graph and a problematic method for edit-distance based similarity search problem. In [41], it additionally gift an even index structure, particularly u – tree, to facilitate effective pruning and economical question process.

K. Gouda et.al [4], proposed a method for calculating graph comparison based on edit distance of the graph. Existing edit distance of a graph calculation methods adopt the best first search or the breadth first search algorithm A Star. These algorithms are space and time bound. At most, these algorithms can compute or calculate the edit distance of the graphs containing twelve vertices in practice. To calculate graph edit similarity calculation on bigger and distant/distinct graphs, [27] presented a method, a novel mapping method which maps the edges. It is used for the computation of graph edit distance. It uses a sub structure isomorphism as a common enumeration solution. In [30], the algorithm uses a backtracking search which is joined with a various number of heuristics/algorithms to reduce space requirements and quickly remove away a big portion of the search space which is being mapped.

X. Zhao et.al [12] proposed a method using Edit Distance for similarity search in graph database. The GED (Graph Edit Distance) is the important function unit in this method. In this method we divide graph data into variable size non overlapping partitions. This is the first method to use variable size partitions of the graph database for graph indexing and filtering. This method also suggest partition based filtering algorithm for reducing time complexity. In GED subgraph search is to retrieve the data graphs that approximately contain

the query graph. The similarity is defined as number edges need to be modified in the subgraphs of graph database respect to the given query graph. This method provide a dynamic partitioning technique to reduce the cost. The cost aware graph partitioning method is the advantage of this method. The updates on the graph database can be handled. It is the advantage of representing data in the graph structure.

L. Hong et.al [6] discussed the query of sub graph matching with set similarity (SMS2) over an oversized graph information, that retrieves sub graphs which are isomorphous structurally to the graph query. In order to do the efficiency method for query of sub graph matching with set similarity question, [36] styled a unique lattice-based index for knowledge graph, and light-weight signatures for each vertices of the query and knowledge vertices. [33] additionally proposed associate degree economical two-phase pruning strategy as well as set similarity pruning as well as pruning in an structural manner.

Zheng [35] proposed a method which is like searching RDF (Resource Description Framework). The RDF data retrieval is used in semantic web like Wikipedia. This method suggested data retrieval in a sub graph matching way. The large knowledge data base is already exist. It will built into a graph data base like converting from RDF data base to graph database. The same way the knowledge data base will be converted into graph database. Because it is easy to query from the graph structured data base comparing to normal knowledge database. The query is converted into graph. In this method after converting the knowledge database to graph database all the subgraphs in graph database are extracted. Compare these subgraphs to the query graph given by the user. In the comparison if result is match then output the matching subgraph.

2.4 GRAPH MINING

S. Zhang et.al[10] proposed that Subgraph query has become most important task in graph matching. Graphs are large and has tens and thousands of vertices and millions of edges. The subgraph indexing and approximate

matching suggest method for indexing subgraph over a very large graph database. Subgraph indexing is to find the query graph in the graph database. It is not necessary to find the exact matching subgraphs from the graph database. This method suggest approximate matching to the query graph. The major difficult task is to index the graph database so that the query matching retrieval from graph database will take less time. For approximate matching to the query graph the Edge Edit Distance is followed. The query graph edge can be modified so that it can match with the subgraphs in the graph database. The Edge Edit should be minimum. So that the user can give the threshold value which is maximum no of edges that can be modified to find the approximate matching subgraphs in the graph database.

2.5 BIOLOGICAL GRAPHS

K. Aoyama et.al [19] discussed about the protein sequence and their enormous development that challenged the bio scientists so that reliable and fully automated methods to quantify huge amount of data sequence in which large amount of work and efforts have been made. Since the practical methods suggested for protein analysis can be doubtful and workload consuming and may be very time consuming, and hence they are not used for identifying efficiently the gene assortment events in a quick span of time. Due to the recent methodologies restrictions, it is a huge task for the people in the scientific crowd to predict and follow the protein analysis and their experimental results. The phylogeny which are obtained are consistently in par with the tree phylogeny which are built by MEGA4 software.

Assayony.M et.al [20] discussed about graph based dividing methods where an vital and static algorithm in computer technology can be used to manipulate the group a large amount of protein sequences or protein sequences that shared a common attribute like similar amino acid sequences. As the size of the graph databases increases largely, a fast and quick graph based protein sequencing grouping method is very much needed to be developed. We have a

parallel and distributed algorithm/procedure in graph-based grouping methods by increasing the performance of available method Protclust [32] by using parallel and distributed methods. The speed of the existing algorithm is being improved, involves using multiple processors without the state of idleness. The complexity of calculating the pair-wise similar is quadratic to the number the sequences compared.

3. PROPOSED SYSTEM DESIGN

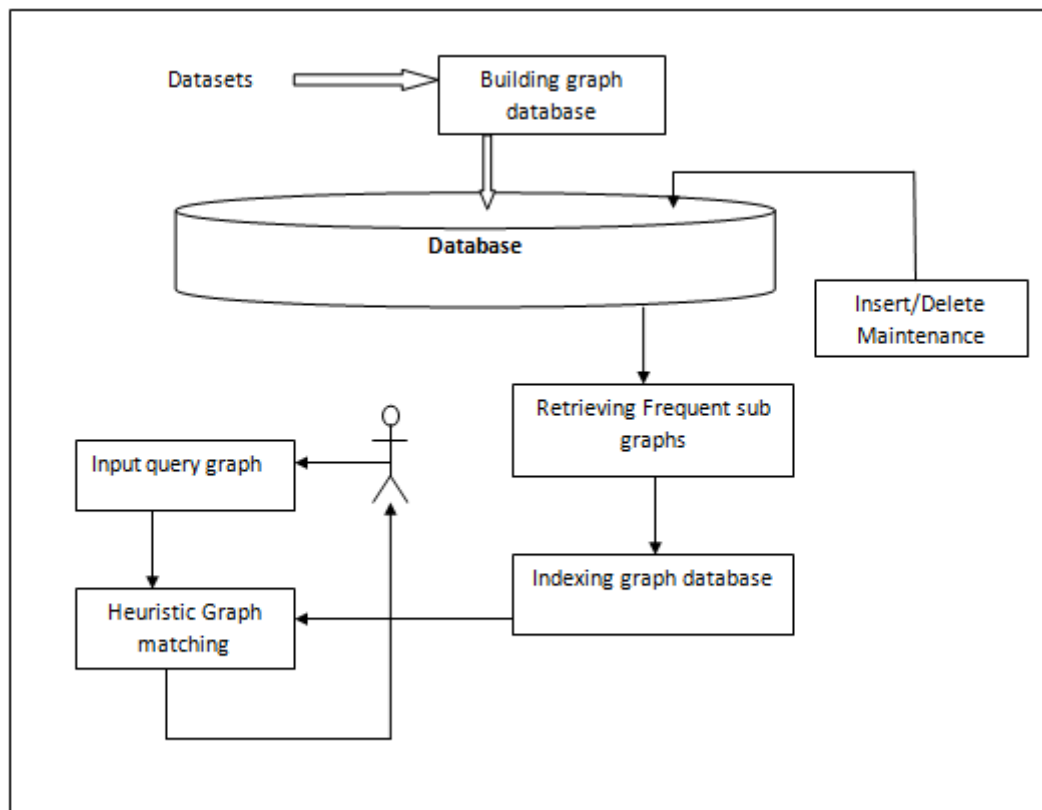


Fig 3.1 Architecture diagram

Datasets are in LINE data format, represented by a list of vertices and edges which is the one of the standard form of representation of graphs. Large graph database is built from the input datasets as shown in Fig 3.1. Then the large graph database is indexed for effective access and retrieval. The graph is further optimized by pruning the dissimilar structures. Using Heuristic Graph

Matching algorithm, best similar graph to the input query graph will be retrieved efficiently. Upon addition of new graphs also, effective insert/delete maintenance is ensured by using the indexes built. The various modules involved are

3.1 Generation of DFS code and retrieving frequent fragments

3.2 Indexing of graph database

3.3 Heuristic graph matching

3.4 Database Maintenance

3.1 GENERATION OF DFS CODE AND RETRIEVING FREQUENT FRAGMENTS

DFS coding is a set of vertices which is used for identifying the graphs. DFS coding is used as the input for generating the frequent sub graphs. DFS coding $(x, y, L_x, L_y, L_{(x, y)})$. The above DFS coding has five parts which are the vertex indexes i and j , the label of the vertex i , the label of the edge and the label of the vertex j . Pre - processing of database has to be done for efficient process before indexing the large graph database. The frequent sub graphs algorithm uses the DFS coding as the input. As shown in Fig 3.2, Frequent sub graph are generated by scanning the Graph database and finding all the edges that can be right most extended to frequent. Then sorting all the frequent set in DFS lexicographical order. This process is repeated for all the tuples identified in the frequent set. The output will contain all the fragments which occurs frequently. The fragments are represented in the DFS coding format in order to understand it easily.

```

Input:  $c$ =DFS code,  $G$ =graph database,  $s$ =min support
Output:  $F$ =frequent sub-graph set;

Algorithm :  $frequent\_subgraph(G, s, c, F)$ 

Scan  $G$  once and find every edge  $e$  so that  $c$  can be right-most extended

If  $c$  can't be extended return

Save  $c \cup e$  into  $c$ ;

Sort  $c$  in DFS lexicographical order

For each tuple in  $c$ 
     $frequent\_subgraph(G, s, c \cup e, F)$ ;
return  $F$ ;

```

Fig 3.2 Frequent sub graph algorithm

3.2 INDEXING OF GRAPH DATABASE

Indexing of a graph database has to be done for efficient retrieval. In this module, the graph database is indexed using the DFS code and the frequent fragments generated. As shown in Fig 3.3, a prefix tree is built using both DFS code and frequent fragments. Prefix tree allows user to check whether a fragment or a sub-graph of a query graph is present in the graph database efficiently. Each node of a prefix tree contains a fragment which may be either frequent or redundant. The prefix tree is built in a manner such that each level contains the fragment of corresponding size. For example, level 1 in the prefix tree contains all the fragments whose size is 1. Level 2 of the prefix tree contains all the fragment which is of size 2 and the fragment should contain the previous level node as a prefix.

Algorithm : Indexing the database

Input: D-Graph database, F-Frequent Fragment

Output: Prefix Tree P

Algo: Indexing (F)

Initialize P with dummy node as root

For each dfscode in set F

P.Insert (dfscode, 0)

Input: dfscode - Array representing a graph

Algo: Insert (dfscode, position)

If position == len (dfscode) return

If dfscode [position] does not exists in P

Add current dfscode in P

Insert (dfscode, position + 1)

Fig 3.3 Indexing algorithm

3.3 HEURISTIC GRAPH MATCHING

Dissimilar structures can be pruned before graph matching. It can be done by pruning the dissimilar structures and selecting the optimal sub graph from various possible graphs. In this module, prefix tree is constructed by inserting the frequent sub-graphs retrieved. Prefix tree is mainly for identifying a super graph in which two given graphs are sub-graphs. After pruning, optimal sub graph can be obtained which can be used for heuristic graph matching. In order to perform pruning efficiently, we can eliminate certain fragments i.e. the main strategy is if a fragment is not present in a prefix tree, we can efficiently avoid all the super graph of the fragment.

After pruning, the input query graph will be given by the user. It uses the frequent sub graphs retrieved from the first module and the input query graph as the input. The optimal sub graph obtained after indexing is compared efficiently with the input query graph and the output gives the resultant graphs which are similar to the query graph upto a given threshold. Fragments of all length up to a given maximum length is compared with the frequent sub-graphs retrieved

and efficient graph matching is done. Graph matching can be done efficiently by eliminating (pruning) the redundant fragments i.e. if a fragment is not in the prefix tree, then its super-graph need not be compared. It can be done by using a hash table which contains the hash codes of the nodes in the prefix tree including the intermediate nodes. Whenever we find a fragment in the query whose hash code does not appear in the hash table, we need not check its super-graphs. As shown in Fig 3.4, fragments of all length up to a given maximum length is compared with the frequent sub-graphs retrieved and efficient graph matching is done.

Algorithm: Graph Matching

Input : D- graph database, q-query graph

Output: Result set R

Let $R = \{\}$

For each fragment x where x is a sub graph of q and $len(x) \leq maxlen$

If x is present in Prefix Tree then

$R = R \cup x$

return R

Fig 3.4 Heuristic Graph Matching Algorithm.

3.4 DATABASE MAINTENANCE

Insertion and Deletion of new incoming graphs has to be done for incremental and updating graphs in the graph database. In this module when the new graphs are added or deleted to the existing graph database it should be updated to the Prefix tree for future indexing and accurate results.

As shown in Fig 3.5, inserting of new graphs takes frequent set as input. A new id is assigned to the incoming graph. Then we compare the fragments to the fragments already in the built Prefix tree and if they match, update the ids of the matching fragment in the id list. In the same way deletion of graphs is also implemented..

Algorithm : Database Maintenance

*Input: Graph Database D, Frequent set F,
Insert(delete) graph g and its id gid, Maximum
fragment size maxL*

*for each fragment x where x is a sub graph of g and
len(x) ≤ maxL do*

If $x \in F$ then

Insert:

Insert gid into the id list of x

Delete:

Delete gid from the id list of x

return;

Fig 3.5 Database Maintenance

4. RESULT AND ANALYSIS

The proposed system of finding the similar graphs from the graph database has better runtime complexity than the existing method. The proposed structure based indexing is evaluated for its performance against the path based indexing and the results are proven to be accurate.

4.1 COMPLEXITY ANALYSIS

The time complexity of various modules are analysed.

4.1.1 Building the prefix tree

The entire frequent fragment set has to be inserted into the prefix tree. So inserting one graph into the prefix tree takes $O(E)$ where E is the number of edges in the graph. Since we have N graphs in the frequent set it requires $O(N * E)$ time to insert all the frequent sub graphs. E is the number of edges in a graph and N is the number of graphs present in the frequent set.

4.1.2 Searching

For generating all the sub graphs of a graph, the worst case time complexity occurs for the graph which is complete. The algorithm does a depth first search for the forward edge alone from all the vertices and it leads to the worst case $O(2^N)$ when the graph is complete. For Non complete graphs it has the best time complexity of $O(V * (V + E))$ which is $O(V^2)$ where V is the number of vertices in the graph and E is the number of edges in the graph.

For searching/querying, a graph will be represented in the form of DFS code which will be given as input. The algorithm needs to traverse all the edges and check whether the traversed prefix is present in the prefix Tree which leads to $O(E)$ where E is the number of edges in the graph.

4.1.3 Retrieving the frequent sub graphs

For retrieving the frequent sub graphs, initially we have to generate all the one edge fragments. For generating it we have to compare all the edges of the graph and sort it. Sorting the edges in the graph takes $O(E \log E)$. Since we have N graphs in the database it leads to $O(N * E \log E)$. Once the single edge fragments are generated it has to be extended. For extending it we need to do a BFS and find the right most edge. It takes $O(V + E)$. Once the vertex is found we have to do a comparison whether this edge can be extended. It leads to $O(V^2)$. It has to be done for all the edges for which we do BFS. If X is the number of one edge fragment generated then it requires $O(X * (V+E) * V^2)$ which is $O(X * V^3)$ where V is the number of vertices in the graph.

4.1.4 Database Maintenance

For inserting a new graph into the database, we have to add the id of the new graph into the prefix tree. So searching has to be done to check whether the graph is present in the prefix tree. It can be done in $O(E)$ where E is the number of edges in the graph.

Number of graphs	Existing(Seconds)	Proposed (Seconds)
5	0.6947	0.5021
10	1.2350	1.0874
15	1.8821	1.5326
20	2.4672	2.0981
25	5.0209	3.4563
30	8.0543	5.3913
35	11.0013	8.9073
40	13.0163	10.2561
45	15.0821	13.4396
50	17.0984	15.2378
55	20.0372	17.4468
60	22.0834	19.5090
70	28.0382	24.4355
80	33.3742	29.9313
90	40.2645	34.0438
100	49.6834	42.3602
140	82.4535	69.2338
180	126.6542	96.3739
200	161.1463	125.8464
225	245.6342	187.2533
250	324.5275	237.4534

Table 4.1 Complexity Analysis – Runtime Comparison

Table 4.1 shows the runtime comparison of path based(Existing) and structure based(Proposed) indexing methods. All the values are compared in seconds.

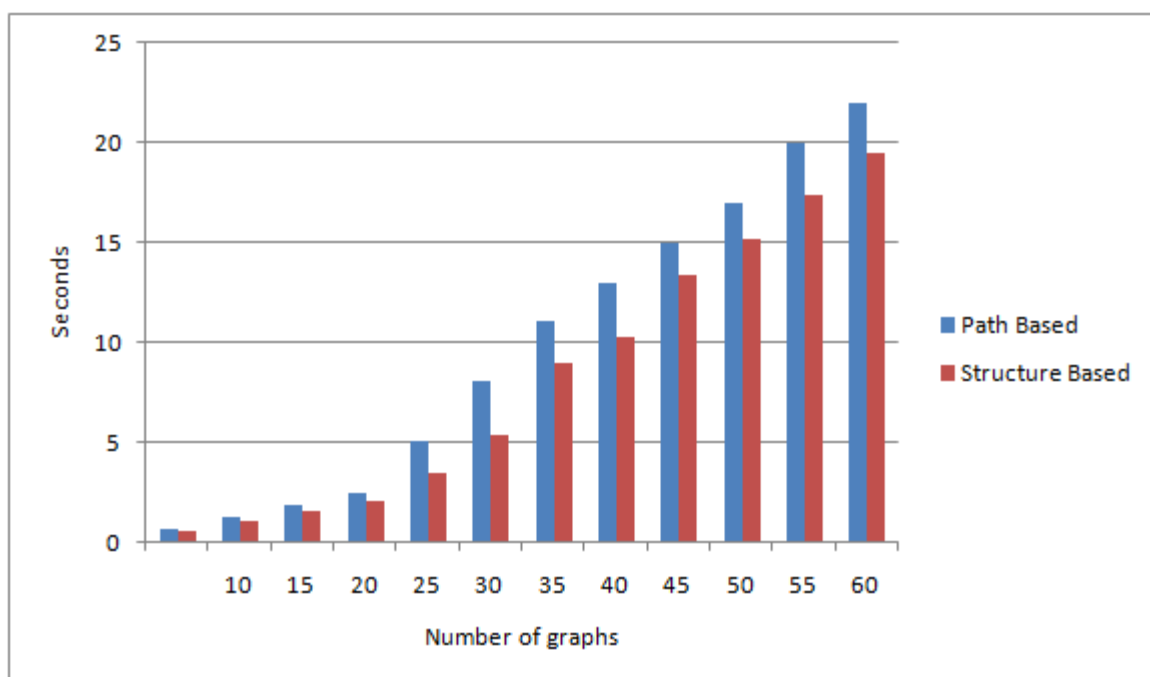


Fig 4.1 Graph Matching – Complexity Analysis

Fig 4.1 shows the graphical representation of path based and structure based indexing methods. All the values are compared in seconds.

Dataset	No. of graphs	Average size		Largest size	
		No. of nodes	No. of edges	No. of nodes	No. of edges
Compound	422	40	42	189	196
Chemical	340	27	28	214	214

Table 4.2 Characteristics of experimental dataset

Table 4.2 shows the number of graphs, largest and average size of the number of nodes and number of edges of the datasets which are taken as samples.

Support Value	Number of frequent fragments generated	
	Dataset 1	Dataset 2
25	693	540
50	517	486
75	487	427
100	415	374
150	298	245
200	234	179
250	167	133
300	106	98
350	77	54
400	48	26

Table 4.3 Comparison of number of frequent fragments generated for different support value

Once the frequent fragments are generated, various query graphs are given and searched with the frequent fragments generated. Table 4.3 also shows that the number of frequent fragment generated for the dataset 1(compound) is larger than the dataset 2(chemical). Various query graphs are matched with the generated fragments and the result shows that the dataset 1 has 83% similarity between the graphs. The dataset 2 has nearly 71% of similarity between its graphs. Similarity is computed by giving half of the dataset as the query graph and other half is considered as a database.

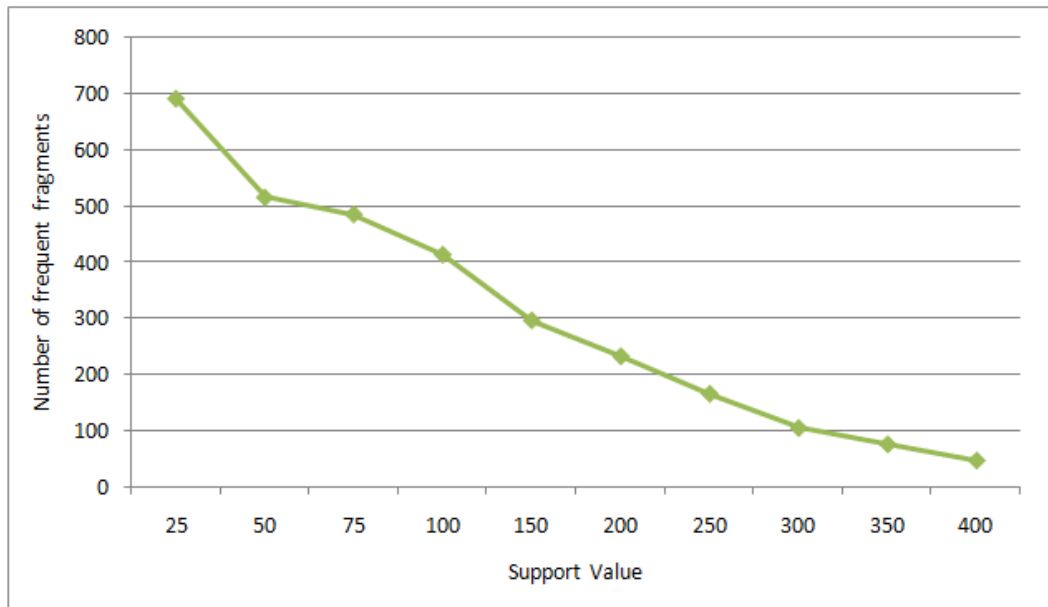


Fig 4.2 Number of frequent fragments generated for different support value

Fig 4.2 shows the number of frequent fragments generated based on the support value. It is clear from the graph that if the support value is inversely proportional to the number of frequent fragments generated i.e if support value is less it generates more number of frequent fragments and vice versa.

5. CONCLUSION AND FUTURE WORK

Graph similarity search is one of the important problems in recent times. It can be applied for medical applications like finding the similar disease patterns. Existing classical algorithms for graph similarity are NP complete. Our method of finding the similar graphs using structure based indexing has a better runtime complexity and high performance compared to the existing methods. Better performance is achieved by using efficient indexing and pruning strategy. Machine learning can also be applied for graph similarity finding. A machine learning algorithm can be used so that we can automate the process of identifying the similar patterns.

References

- [1] W. Zheng, L. Zou, X. Lian, D. Wang and D. Zhao, "Efficient Graph Similarity Search Over Large Graph Databases," in IEEE Transactions on Knowledge and Data Engineering, vol. 27, no. 4, pp. 964-978, April 1 2015.
- [2] E. S. S. Dongoran, W. K. Rahmat Saleh and A. A. Gozali, "Analysis and implementation of graph indexing for graph database using GraphGrep algorithm," 3rd International Conference on Information and Communication Technology (ICoICT), pp. 59-64, 2015.
- [3] Xifeng Yan Philip S. Yu, Jiawei Han " Graph Indexing: A Frequent Structure-based Approach," IEEE International Conference on Big Data, vol., no.33, pp.271,280, 2014.
- [4] K. Gouda and M. Hassaan, "CSI_GED: An efficient approach for graph edit similarity computation," 32nd International Conference on Data Engineering (ICDE), Helsinki, pp. 265-276, 2016.
- [5] G. Wang, B. Wang, X. Yang and G. Yu, "Efficiently Indexing Large Sparse Graphs for Similarity Search," in IEEE Transactions on Knowledge and Data Engineering, vol.24, no.3, pp.440-451, March 2012.
- [6] L. Hong, L. Zou, X. Lian and P. S. Yu, "Subgraph Matching with Set Similarity in a Large Graph Database," in IEEE Transactions on Knowledge and Data Engineering, vol. 27, no. 9, pp. 2507-2521, Sept1 2015.
- [7] Patil, Shefali, Vaswani, Gaurav; Bhatia, Anuradha , "Graph Databases- An Overview" in International Journal of Computer Science & Information Technology ,Vol. 5 Issue 1, p657, 2014.

- [8] Xiaoli Wang, Xiaofeng Ding, Anthony K.H. Tung, "An Efficient Graph Indexing Method", School of Computer Science, Huazhong University of Science and Technology, P. R. China, pp. 456-901, 2010.
- [9] L.Zou, J.Mo, L.Chen, M.T. Ozsü and D. Zhao, "gStore: Answering SPARQL queries via subgraph matching," Proc. VLDB Endowment, vol. 4, no. 8, pp. 482–493, 2011.
- [10] S. Zhang, J.Yang, and W. Jin, "Sapper: Subgraph indexing and approximate matching in large graphs," Proc. VLDB Endowment, vol. 3, no. 1, pp. 1185–1194, 2010.
- [11] X. Zhao, C.Xiao, X. Lin, and W. Wang, "Efficient graph similarity joins with edit distance constraints," in Proc.IEEE 28th Int. Conf. Data Eng, pp. 834–845, 2012.
- [12] X. Zhao, C.Xiao, X. Lin, Q. Liu, and W. Zhang, "A partition-based approach to structure similarity search," Proc. VLDB Endowment, vol. 7, no. 3, pp. 169–180, 2013.
- [13] X. Wang, X. Ding, A. K. H. Tung, S. Ying, and H. Jin, "An efficient graph indexing method," in Proc. IEEE 28th Int. Conf. Data Eng, pp. 210–221, 2012.
- [14] Dong Dai, Yong Chen, Kimpe, D. Ross, "Provenance-based object storage prediction scheme for scientific big data applications," IEEE International Conference on Big Data, vol no. 4, pp.271,280, 2014.

- [15] Williams, D.W. Jun Huan Wei Wang, "Graph Database Indexing Using Structured Graph Decomposition", IEEE Transactions on Data Engineering, pp.976-985, 2012.
- [16] K. McGarry and U. Daniel, "Computational techniques for identifying networks of interrelated diseases," 14th UK Workshop on Computational Intelligence, Bradford, pp. 1-8, 2014.
- [17] Hung, Benjamin WK, Anura P. Jayasumana, and Vidarshana W. Bandara. "Pattern Matching Trajectories for Investigative Graph Searches." IEEE International Conference on Data Science and Advanced Analytics (DSAA), pp. 71-79, IEEE, 2016.
- [18] G. Wang, B. Wang, X. Yang and G. Yu, "Efficiently Indexing Large Sparse Graphs for Similarity Search," in IEEE Transactions on Knowledge and Data Engineering, vol. 24, no. 3, pp. 440-451, 2013.
- [19] K. Aoyama, S. Watanabe, H. Sawada, Y. Minami, N. Ueda and K. Saito, "Fast similarity search on a large speech data set with neighbourhood graph indexing," IEEE International Conference on Acoustics, Speech and Signal Processing, Dallas, TX, pp. 5358-5361, 2010.
- [20] Assayony, M.; Rashid, N.A., "Design of a parallel graph-based protein sequence clustering algorithm", IEEE Transactions on Information Technology, vol.3, pp.1-8, 26-28, 2012.
- [21] J. Rocha, "Graph Comparison by Log-Odds Score Matrices with Application to Protein Topology Analysis," in IEEE/ACM Transactions on Computational Biology and Bioinformatics, vol. 8, no. 2, pp. 564-569, 2011.

- [22] H. Li and C. Liu, "Prediction of protein structures using a map-reduce Hadoop framework based simulated annealing algorithm," IEEE International Conference on Bioinformatics and Biomedicine, Shanghai, pp. 6-10, 2013.
- [23] DePiero, F.W.; Carlin, J.K., "Structural Matching Via Optimal Basis Graphs", IEEE Transactions on Pattern Recognition, vol.3, pp.449-452, 2015.
- [24] F. Bai, Y. Li and H. Gao, "The Comparison of Protein Secondary Structure Based on the Graphical Representation," International Conference on Computational and Information Sciences, Chengdu, China, pp. 11-14, 2011.
- [25] Changjun Wu; Kalyanaraman, A.; Cannon, W.R., "pGraph: Efficient Parallel Construction of Large-Scale Protein Sequence Homology Graphs", IEEE Transactions on Parallel and Distributed Systems, vol.23, no.10, pp.1923-1933, 2012.
- [26] H. Hu; Z. Li; H. Dong; T. Zhou, "Graphical Representation and Similarity Analysis of Protein Sequences Based on Fractal Interpolation," in IEEE/ACM Transactions on Computational Biology and Bioinformatics, pp. 99-111, 2012.
- [27] Jiefeng Cheng; Yu, J.X.; Bolin Ding; Yu, P.S.; Haixun Wang, "Fast Graph Pattern Matching", IEEE Transactions on Data Engineering, pp.913-922, 2013.
- [28] J. Wu and DaiChuan Ma, "Comparisons of the protein-protein interaction networks constructed from the DIP database with different version," International Conference on Electric Information and Control Engineerin, Wuhan, pp. 236-239, 2011.
- [29] Li Chen; Gupta, A.; Kurul, M.E., "Efficient algorithms for pattern matching on directed acyclic graphs" IEEE Transactions on Data Engineering, pp.384-385, 2015.

[30] M. Mernberger, G. Klebe and E. Hullermeier, "SEGA: Semiglobal Graph Alignment for Structure-Based Protein Comparison," in IEEE/ACM Transactions on Computational Biology and Bioinformatics, vol. 8, no. 5, pp. 1330-1343, 2011.

[31] P. Li, L. Heo, M. Li, K. H. Ryu and G. Pok, "Protein function prediction using frequent patterns in protein-protein interaction networks," Eighth International Conference on Electric Information and Control Engineerin, Shanghai, pp. 1616-1620, 2011.

[32] Yoo, Young Joon, TusharSandhan, Jinyoung Choi, and Sun Kim. "Towards simultaneous clustering and motif-modeling for a large number of protein family." IEEE International Conference on Bioinformatics and Biomedicine (BIBM), pp. 22-28, 2013.

[33] Hu, Hailong, Zhong Li, Hongwei Dong, and Tianhe Zhou. "Graphical Representation and Similarity Analysis of Protein Sequences Based on Fractal Interpolation." IEEE/ACM Transactions on Computational Biology and Bioinformatics , vol no. 1 (2017): 182-192, 2017.

[34] Sheng-Lung Peng, Yu-Wei Tsay, "Using Bipartite Matching in Graph Spectra for Protein Structural Similarity", IEEE Transactions on BiomedicalEngineering and Informatics, pp.495-500, 16-18, 2013.

[35] Zheng, Weiguo, Xiang Lian, Lei Zou, Liang Hong, and Dongyan Zhao. "Online Subgraph Skyline Analysis Over Knowledge Graphs." IEEE Transactions on Knowledge and Data Engineering, vol no. 7, pp 1805-1819, 2016.

- [36] Thomas, Minta, AnneleenDaemen, and Bart De Moor. "Maximum likelihood estimation of GEVD: Applications in Bioinformatics." *IEEE/ACM Transactions on Computational Biology and Bioinformatics* 11, vol no. 4, pp. 673-680, 2014.
- [37] Mernberger, Marco, Gerhard Klebe, and EykeHullermeier. "SEGA: Semiglobal Graph Alignment for Structure-Based Protein Comparison." *IEEE/ACM transactions on computational biology and bioinformatics* 8, vol no.5 , pp. 1330-1343, 2011.
- [38] Wu, Wei, Anuj Srivastava, Jose Laborde, and Jinfeng Zhang. "An efficient multiple protein structure comparison method and its application to structure clustering and outlier detection." *IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pp. 69-73. IEEE, 2013.
- [39] Li, Peipei, LyongHeo, Meijing Li, Keun Ho Ryu, and GoucholPok. "Protein function prediction using frequent patterns in protein-protein interaction networks.", *Eighth International Conference on Fuzzy Systems and Knowledge Discovery (FSKD)*, vol. 3, pp. 1616-1620, IEEE, 2011.
- [40] Qiao, Rui, XiaoleiZhong, Ling Zhang, and Heng He. "Graph Pattern Matching through Model Checking.", *8th International Conference on Database Theory and Application (DTA)*, pp. 1-5. IEEE, 2015.
- [41] Fairey, Jason, and Lawrence Holder. "StarIso: Graph Isomorphism Through Lossy Compression." *Data Compression Conference (DCC)*, pp. 589-589. IEEE, 2016.
- [42] Zhong, Haojian, Lida Xu, Cheng Xie, Boyi Xu, Fenglin Bu, and HongmingCai. "A Similarity Graph Matching Approach for Instance Disambiguation.", *4th International Conference on Enterprise Systems (ES)*, pp. 21-28. IEEE, 2016.

