

# Optimizing $K^2$ Trees: A Case for Validating the Maturity of Network of Practices

Quan Shi<sup>1</sup>, Yanghua Xiao<sup>2\*</sup>, Bessis Nik<sup>3\*</sup>, Yiqi Lu<sup>2</sup>, Yaoliang Chen<sup>2</sup> and Hill Richard<sup>3</sup>

<sup>1</sup> School of Computer Science and Technology, Nantong University, Nantong China

<sup>2</sup> School of Computer Science, Fudan University, Shanghai, China

<sup>3</sup> School of Computing and Mathematics, University of Derby, United Kingdom

[sq@ntu.edu.cn](mailto:sq@ntu.edu.cn); [shawyh@fudan.edu.cn](mailto:shawyh@fudan.edu.cn); [n.bessis@derby.ac.uk](mailto:n.bessis@derby.ac.uk);

[luyiqi@gmail.com](mailto:luyiqi@gmail.com); [chenyao2@gmail.com](mailto:chenyao2@gmail.com); [r.hill@derby.ac.uk](mailto:r.hill@derby.ac.uk)

## Abstract

Of late there has been considerable interest in the efficient and effective storage of large-scale network graphs, such as those within the domains of social networks, web and virtual communities. The representation of these data graphs is a complex and challenging task and arises as a result of the inherent structural and dynamic properties of a community network, whereby naturally occurring churn can severely affect the ability to optimize the network structure. Since the organization of the network will change over time, we consider how an established method for storing large data graphs ( $K^2$  tree) can be augmented and then utilized as an indicator of the relative maturity of a community network. Within this context, we present an algorithm and a series of experimental results upon both real and simulated networks, illustrating that the compression effectiveness reduces as the community network structure becomes more dynamic. It is for this reason we highlight a notable opportunity to explore the relevance between the  $K^2$  tree optimization factor with the maturity level of the network community concerned.

**Keywords**  $K^2$  tree; storage optimization; DFS code; compression algorithm; network of practices

## 1. Introduction

The efficient organization and representation of graph data and the corresponding data processing method is a topical issue for the database research community. In recent years, the exponential rate of growth and pervasiveness of the Internet has resulted in larger scale data sets, particularly since the advent of social networking applications. For instance, Facebook has in excess of 484 million registered users, with each individual having, on average, 120 ‘friends’ [1].

Clearly the volume and complexity of these datasets has much to offer the research community, in terms of understanding established and emerging developments in data organization and usage. In contrast with more traditional repositories such as geographical map stores, these large-scale social networks have the potential to exhibit an increased richness that has not been witnessed before. Existing research can be divided into two categories. In the first category, the methods focus upon graph representation on secondary storage devices. For example, Aggarwal et al. [2] build a connectivity index for massive disk-resident graphs based on a digest graph, which can be directly loaded into memory. The second category considers graph compression as a means of reducing storage costs. Both of these approaches tend to explore the property of an adjacency list

---

\* Corresponding author: Dr Yanghua Xiao, School of Computer Science, Fudan University, Email: [shawyh@fudan.edu.cn](mailto:shawyh@fudan.edu.cn); Prof Nik Bessis, University of Derby, School of Computing & Maths, Kedleston Road, Derby, DE22 1GB, United Kingdom, Email: [n.bessis@derby.ac.uk](mailto:n.bessis@derby.ac.uk), Tel.: 00441332592108, Fax: 00441332597741

or matrix. Amongst them, representative methods include:  $K^2$  tree [3]; compression based upon vertex similarity [4];  $\zeta$ -code [5]; use of frequent patterns [6]; compression using graph skeleton by symmetry [7]; social network compression using shingle ordering [8]; and finally a back-links schema and hybrid compression method utilizing dynamic arrays [9]. The next section will offer an introduction to the contextual basis of our work.

## **2. Contextual basis**

### **2.1. Networks and communities**

There has been a considerable amount of research with regard to the concept of Communities of Practice (CoP). This work takes the perspective that a CoP has a social learning theory dimension [10] [11] [12] which in effect acknowledges elements of knowledge transformation between connected nodes on the basis of their shared interest.

The emergence of Web technologies has concentrated the study of CoP from various technological based developments, including web and virtual communities, online communities and social networks. For brevity we shall use the term 'Network of Practice' (NoP) [13] to refer to the overall set of informal, emergent social networks that facilitate learning and knowledge sharing between individuals conducting practice-related tasks. [12] describes the life-cycle of a CoP as an informal learning and development process encompassing three key stages. We will use the three key stages to annotate user activity and participation within any network of practice (NoP).

During the first stage, a NoP is in its infancy and comprises a scattered set of individuals across a particular network space, that have intent to gain awareness of their space. This is characterized by behavior that promotes and accelerates the process of familiarization between different parties, for instance by flooding the community network space with message requests for information.

Following on and as a result of stage one interaction, an individual and collective awareness emerges that facilitates an active environment. Individual networked nodes will exhibit more bi-directional communication, with less emphasis upon the broadcasting of messages to the masses. This is where the network space exhibits evidence of user participation and thus, individuals working towards the fulfillment of their goals.

Finally, the individuals demonstrate more fruitful communicative relationships, and commonality between clusters of activity becomes more apparent suggesting the emergence of a particular expertise. The predominant behavior in the network at this stage is that of relationship maintenance; rather than relationship propagation, there is little need to broadcast requests.

Lately, there is an increased interest in understanding what constitutes 'success' in a NoP. Most studies describe success of communities from the perspective of the information system success models. [10] points out that the main concern with this model is that it does not consider the social relationships among members nor the structure of the networked community. For this reason, [10] [11] suggest a new approach based on social network analysis (SNA). Following their work, we are particularly interested in understanding the maturity of a NoP which in effect reflects its relative success.

Specifically, the ability to make an assessment of the maturity of a CoP presents an opportunity for emerging systems to develop more sophisticated behaviors, such as self-awareness. Within this context, we are interested in the development of an optimized  $K^2$  tree method that can be utilized to demonstrate the maturity of a NoP. The remainder of this article describes how we have augmented the simple  $K^2$  tree approach in order that more complex, community network graphs can be compressed effectively for storage, as well as providing an indication as to the relative

maturity of that network. First, we introduce the  $K^2$  approach.

## 2.2. $K^2$ tree approach

If a real-world network is represented as a data graph, it will typically demonstrate empty regions. A  $K^2$  tree compresses a large number of these zero value regions in an adjacency matrix, in order that a reduced number of  $K^2$  tree nodes will result, thus demonstrating effective compression.

However, a simple  $K^2$  tree is still too abstract for the effective representation of the structural characteristics of real-world networks. In terms of compression efficiency, there is still considerable potential to improve upon a simple  $K^2$  tree. Firstly, many real networks have an evident hierarchy of communities. Nodes inside one community are highly interconnected with other nodes, while links between different communities are relatively sparse.

If it were possible to sort the nodes of a graph, based upon the similarity of their respective indices, then the corresponding adjacency matrix would reflect a concentration of the potential values, in this case either containing values of 'one' or 'zero'. It follows that this concentration of values within the matrix would significantly improve the resulting compression effectiveness.

However, there are two fundamental issues with this approach. Firstly, a simple  $K^2$  tree cannot represent the complex network structure of a community. Secondly, the established  $K^2$  compression method relies upon a fixed value for  $K$ ; should the structure of the community alter, and therefore the distribution of values in the adjacency matrix change, the calculated value for  $K$  would now be inappropriate resulting in sub-optimal compression.

This presents a challenge that forms the basis of this research, in which we not only explore the optimization of the  $K^2$  method in order maximize graph compression for storage, but we also take account of the need to represent the maturity of more complex, dynamic community network structures such as a NoP.

### 2.2.1. Basic concepts

We consider an undirected graph as an exemplar by which we can discuss the challenges of optimization using a  $K^2$  tree. The proposed algorithm can also be extended to solve the difficulties of representation and storage of graphs with multi-edges and self-loops. This section gives some basic concepts and the formal definitions used in this paper, such as the undirected graph and a  $K^2$  tree.

**Definition 1:** Undirected graph. An undirected graph is a tuple  $G = (V, E)$ , where  $V$  is the set of vertices in graph  $G$  and  $E \subseteq V \times V$  is the edge set.

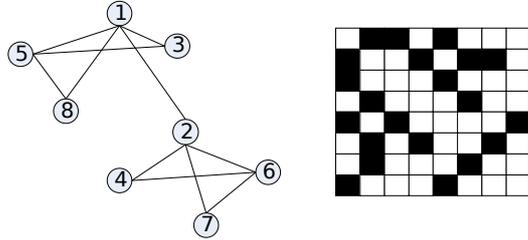
For graph  $G$ , its corresponding adjacency matrix  $A$  is a  $N \times N$  matrix where  $N$  is the vertex number of graph  $G$  and  $a_{ij} = 1$  iff  $(v_i, v_j) \in E$ , otherwise  $a_{ij} = 0$ . Figures 1(a) and 1(b) gives an example of the graph and its corresponding adjacency matrix.

**Definition 2:**  $K^2$  Tree. An adjacency matrix can be represented by an unbalanced  $K^2$ -dimension tree, called a  $K^2$  tree. For a  $K^2$  tree, in addition to the leaf nodes used to represent the adjacency matrix, the other nodes can be represented by 0 for a leaf node and 1 for an internal node. The lowest level in the tree is the root, followed by the second level ( $K^2$  children of root), each of which is either 0 or 1. In a  $K^2$  tree, every internal node has  $K^2$  children.

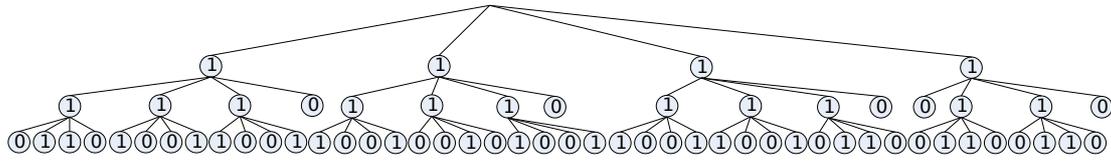
Given an adjacency matrix, the construction of a  $K^2$  tree is similar to that of QuadTree [14]. First the adjacency matrix is divided equally into parts, that is, divide the original matrix into  $K^2$   $K \times K$  sub-matrices. These  $K^2$  sub-matrices correspond to the  $K^2$  children of the root of  $K^2$  tree. If a sub-matrix has a value of 1, the value of the child node is also 1, otherwise it is assigned a 0. It follows

that, a node of value 0 means that the elements in a sub-matrix it represents are all 0, which as a consequence, makes this a leaf node. For those nodes of value 1, we recursively divide the sub-matrix it represents into  $K^2$  parts until all of the values of a sub-matrix are equal to 0, or a sub-matrix which consists of only one element is reached.

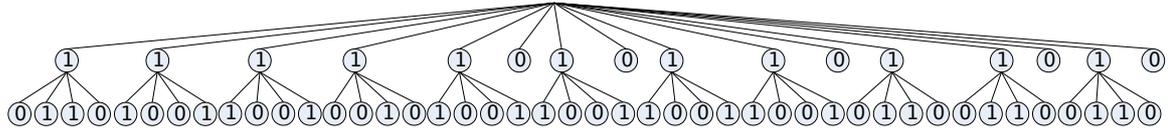
For a  $K^2$  tree, a large  $K$  value leads to a tree with relatively small number of layers but a large spread of leaf nodes. Figures 1(c) and 1(d) show the  $K^2$  tree when the  $K$  value is 2 and 4.



(a) Sample graph and (b) related adjacency matrix



(c)  $K^2$  tree when  $K=2$



(d)  $K^2$  tree when  $K=4$

Fig.1 A sample graph and its  $K^2$  trees

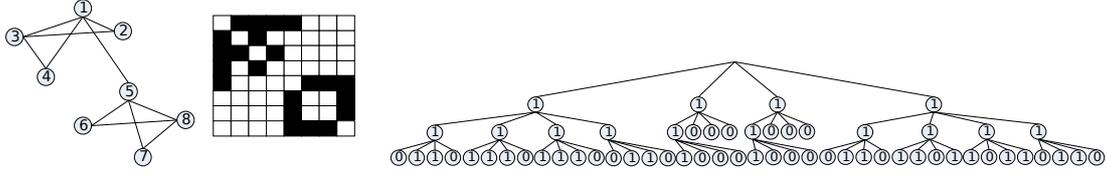
### 3. Optimization of $K^2$ tree

In this section, we present various optimization techniques for  $K^2$  tree.

#### 3.1. DFS Code with heuristic rule

For sparse networks that exhibit a community structure, there must exist some specific node orders which make the elements of value 1 in the corresponding adjacency matrix relatively concentrated in some sub-matrices, rather than randomly appearing in the adjacency matrix. Therefore, by finding an effective node encoding and reordering schema the number of internal nodes in  $K^2$  Tree can be suitably reduced.

For a graph with  $N$  vertices, there are  $N!$  possible node orders in total, each corresponding to a unique  $K^2$  tree. It is clear that a different node order leads to different  $K^2$  tree size. As shown in Figure 2, after reordering the nodes in the graph of Figure 1(a), the  $K^2$  tree ( $K=2$ ) obtained has 4 nodes less than the original. Amongst all possible orders of nodes, the one leading to  $K^2$  tree with minimum tree node number is the optimal node order. Finding the optimal node order is intractable. Even when  $K=2$ , this optimal node ordering problem can be reduced to a minimum bisection problem, which has been proved to be NP-hard [15]. When  $K$  is dynamic, the problem becomes more complex. Therefore we need to use heuristic rules to find the approximate optimal solution. This article uses depth first search (DFS) with heuristic rules to encode the vertices in a graph, and then takes the DFS order as approximate for the optimal node order.



(a) Reordered Graph (b) Corresponding adjacency matrix (c)  $K^2$  tree when  $K = 2$

Fig.2 Reordered Graph and its corresponding  $K^2$  tree

The optimal node order should be used to encode nodes inside a community continuously. Therefore an appropriate DFS heuristic rule to guide the process of traversing the nodes within one community is needed. As such, the index of these nodes in one community is continuous which makes them more likely to be in one sub-matrix when the original matrix is being divided. In our experiment, it is shown that this heuristic rule will achieve an improved performance together with the retention of community structural characteristics. We use a simple and effective heuristic rule: *structural similarity*. For a given node pair  $(u, v)$ , the structural similarity between  $u$  and  $v$  is  $\alpha(u, v)$ ,

$$\alpha(u, v) = \frac{|N(u) \cap N(v)|}{|N(u) \cup N(v)|} \quad \square \square,$$

where  $N(u)$  denotes the neighbor set of node  $u$ . For every edge  $e(u, v)$ , we can define its structural similarity  $\alpha(u, v)$ . Using DFS, the neighbor which has the largest structural similarity is selected as the next node to visit.

We take Figure 1(a) as an example. In Figure 1(a), we have two communities where the indices of nodes are randomly assigned. Figure 2(a) shows the re-encoded graph from node 1 using our heuristic DFS rule. Every node inside a community is indexed continuously. Its corresponding adjacency matrix is given in Figure 2(b). Edges inside one community tend to appear in the same block, which will assist the compression efficiency of  $K^2$  tree. The graph in Figure 2(c) has 4 less nodes than that in graph 1(c).

Structural similarity has been widely used in computing the collection similarity, which is also known as the Jaccard coefficient. For graphs represented by an adjacency list, we can first perform a sort on an adjacency list. For an edge  $(u, v)$ , its Jaccard coefficient can be obtained by using mergesort in  $\Theta(d(u)+d(v))$ . Based on mergesort, it can be proved that we can obtain the Jaccard coefficient of every edge in the graph in  $O(N \log^2 N)$ .

**Lemma 1:** For a sparse graph  $G(M=\Theta(N \log N))$ , where  $N$  is the number of vertices,  $M$  is the number of edges. Based on mergesort, the time complexity of calculating the Jaccard coefficient of all edges is  $O(N \log^2 N)$ .

**Proof:**

The cost of sorting an adjacency list is

$$\sum_{u \in V} d(u) \log d(u) \leq \log d_{\max} \sum_{u \in V} d(u) \leq M \log d_{\max}$$

where  $d_{\max}$  is the maximal degree in the graph. Therefore the total cost of mergesort is:

$$\sum_{(u,v) \in E} d(u) + d(v) = \sum_{u \in V} d(u)^2 \leq \left( \sum_{u \in V} d(u) \right)^2 / N = M^2 / N$$

Because of  $M=\Theta(N \log N)$ , the time complexity is  $O(N \log^2 N)$ .

Large scale, real world networks such as social networks always have a *scale-free network* topology[16]. Its degree distribution obeys a power law form:

$$d(v) = \frac{1}{N^R} (r(v))^R$$

where  $r(v)$  is the degree rank of vertex  $v$  in the graph,  $R$  is a constant and  $R < 0$ . For real graphs of this kind, it can be proved that the Jaccard coefficient can be calculated in  $O(N \log N)$  when  $R \leq -0.5$ . In fact, the  $R$  value of most of the scale-free networks is within  $-2$  to  $-3$  [17].

**Corollary 1**(Time complexity of calculating the Jaccard coefficient on scale-free networks):

For graph with degree distribution of  $d(v) = \frac{1}{N^R} (r(v))^R$  and  $M = \Theta(N \log N)$ , the time complexity of calculating the Jaccard coefficient of all edges based on mergesort is  $O(N \log N)$ .

**Proof:**

The cost of sort on adjacency list is the same as the proof in lemma 1, which is  $O(N \log N)$ . The accumulated cost is:

$$\sum_{(u,v) \in E} d(u) + d(v) = \sum_{u \in V} d(u)^2 = \sum_{i=1}^N i^{2R} / N^{2R}$$

Therefore, when  $R \leq -0.5$  the upper bound is  $O(N \log N)$ .

For large scale networks with more than 1M nodes the method given above to compute the Jaccard coefficient is still too costly. An alternative is to adopt shingle [18] instead. Shingle has been used widely in similarity computation for large scale networks. It [19] has been proved that for a random permutation  $\pi$  on a given set  $U = A \cup B$ ,

$$\Pr[\pi^{-1}(\min_{a \in A} \{\pi(a)\}) = \pi^{-1}(\min_{b \in B} \{\pi(b)\})] = \frac{|A \cap B|}{|A \cup B|} \square 2 \square$$

That is to say, the probability that the minimum element of  $A$  under permutation  $\pi$  equals to the minimum element of  $B$  under permutation  $\pi$  is the Jaccard coefficient of set  $A$  and  $B$ . Based on the theorem above, the  $O(N)$  shingle computation method has been proposed [20] to obtain the Jaccard coefficient. In this case, the time complexity of computing the Jaccard coefficient in our heuristic rule is  $O(N)$ .

### 3.2. Self-adaptive K

Another key factor of  $K^2$  tree compression is the selection of  $K$ , which is usually fixed for a simple  $K^2$  tree. In many cases, a fixed  $K$  may break a sub-matrix which has many values of one into a scatter, resulting in an incremental increase in the number of  $K^2$  tree internal nodes. As shown in Figure 3, if  $K=3$ , we need to divide the matrix according to the community in this graph, otherwise the community will be broken. Therefore an appropriate value for  $K$  is vital.

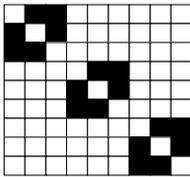


Fig. 3 A 9x9 adjacency matrix

To address this, we consider  $K$  no longer to be fixed value, which leads to a variable tree structure with every internal node having a different fan-out number. Ideally, an algorithm has to propose the best fan-out number for every internal node, thus presenting a significant computational overhead. We propose that only the sibling nodes that share a common parent node will share also

the same value for  $K$ . In this section, we will illustrate by means of theoretical analysis and experimental results that the cost of this is relatively low, whilst also achieving good performance.

---

**Algorithm 1: BuildK<sup>2</sup>Tree( $G$ )**

---

**Input:** Graph  $G$

**Output:**  $TG$ : Optimal K<sup>2</sup> tree;

---

```

1 □ Let  $M$  be the adjacent matrix of  $G$ ;
2 □  $TG\_set = \{\}$ ;
3 □ For each  $k$  from  $k_{min}$  to  $k_{max}$ 
4 □  $TG\_set = TG\_set \cup \{\text{ConstructK}^2(M, k, 0)\}$ ;
5 □ Return the tree with minimal nodes in  $TG\_set$  as  $TG$ ;
```

---

The algorithm for building a K<sup>2</sup> Tree with a self-adaptive K value is described in Algorithm 1. We try every possible k value in  $[k_{min}, k_{max}]$ , then call the method **ConstructK<sup>2</sup>** to calculate the optimal K<sup>2</sup> tree, before selecting the best value for K. In Algorithm 2, lines 1 and 4 iterate through the current sub-matrix to determine whether it is all zeroes or ones, thus identifying a leaf node. Otherwise, go to line 8 to recursively compute the best value for K. Algorithm 2 also attempts to divide the current sub-matrix into every possible value for K, recording the resulting node number (lines 11~13) and storing as *descendant\_count*. After all possible values for K have been tried, Algorithm 2 takes the value stored in *descendant\_count* and stores the minimum value into *m[level]* (line 17~19). This results in an optimized K<sup>2</sup> tree (line 20~21).

---

**Algorithm 2: ConstructK<sup>2</sup>( $M, K, level$ )**

---

**Input:** Adjacency matrix  $M$ ,  $K$ (current K value), *level*(current level)

**Output:**  $TG$ : optimal K<sup>2</sup> tree;

---

```

1 □ IF all the elements in  $M$  are 0 THEN {
2 □  $m[level] = 0$ ; //m[level]store the number of nodes of the
   K2 sub-trees corresponding to  $M$ 
3 □ Construct a leaf with label 0 as TG; }
4 □ IF all the elements in  $M$  are 1 THEN {
5 □  $m[level] = 0$ ;
6 □ Construct a leaf with label 1 as TG; }
7 □ Else
8 □  $min\_TG\_set = \{\}$ ;
9 □  $m[level] = \infty$ ;
10 □ Divide  $M$  into K2 SubMatrixes with equal size.
11 □ For each  $k$  from  $k_{min}$  to  $k_{max}$  {
12 □  $TG\_set = \{\}$ ; //store the K2 tree of submatrix
13 □  $descendant\_count = 0$ ;
14 □ For each SubMatrix {
15 □  $TG\_set = TG\_set \cup \{\text{ConstructK}^2(\text{SubMatrix},$ 
    $k, level+1)\}$  □
16 □  $descendant\_count += m[level + 1] + 1$ ; }
17 □ IF  $m[level] > descendant\_count$  THEN {
18 □  $min\_TG\_set = TG\_set$ ;
19 □  $m[level] = descendant\_count$ ; }
```

---

---

```

20 □ Create a new tree  $TG$  with a root node  $R$ ;
21 □ Set  $R$ 's subtrees as  $Min\_TG\_Set$ ;
22 □ Return  $TG$ ;

```

---

The recurrence formula is given below:

$$T(n) = d_k k_i^2 T\left(\frac{n}{k_i^2}\right) + C \quad \square 3 \square,$$

where  $k_i$  is the current recurrence depth.

Constant  $C$  corresponds to lines 6 and 9 in Algorithm 2. By iterating we can determine whether the sub-matrix contains either all zeroes or all ones. In formula 3,  $k_i$  corresponds to a different recurrence depth, but in practice,  $k$  is selected within a small range. In later experiments, we select  $k$  in  $[2,4]$ . In our analysis, we assume selection of the same  $k_i$  for every recurrence depth such as  $K_{mid}$ , and then according to the master theorem, the complexity of our algorithm is  $O(N^\alpha)$ ,

where  $\alpha = \log_{K_{mid}^2} d_k k_{mid}^2 = 1 + \frac{1}{2} \log_{K_{mid}} d_k$ . In subsequent experiments, we select a  $k$  value in  $[2,4]$ , which provides good performance in our experiment on a network with  $10^5$  nodes where the complexity is  $O(N^{1.5})$ .

## 4. Experiments

We implemented the proposed algorithm using C++. All the experiments were conducted on Windows 7 Professional with Intel Core Duo 1.60GHz and 2GB memory.

### 4.1. Dataset

The network data used in this experiment is shown in Table 1, where  $N$  is the vertex number,  $M$  is the edge number and  $2M/N$  is the average degree. Gq is the random network generated by Pajek and Lesmis[21]. Football[22], Cond\_mat[23] and DBLP[24] is real network data.

Table 1 Experimental data

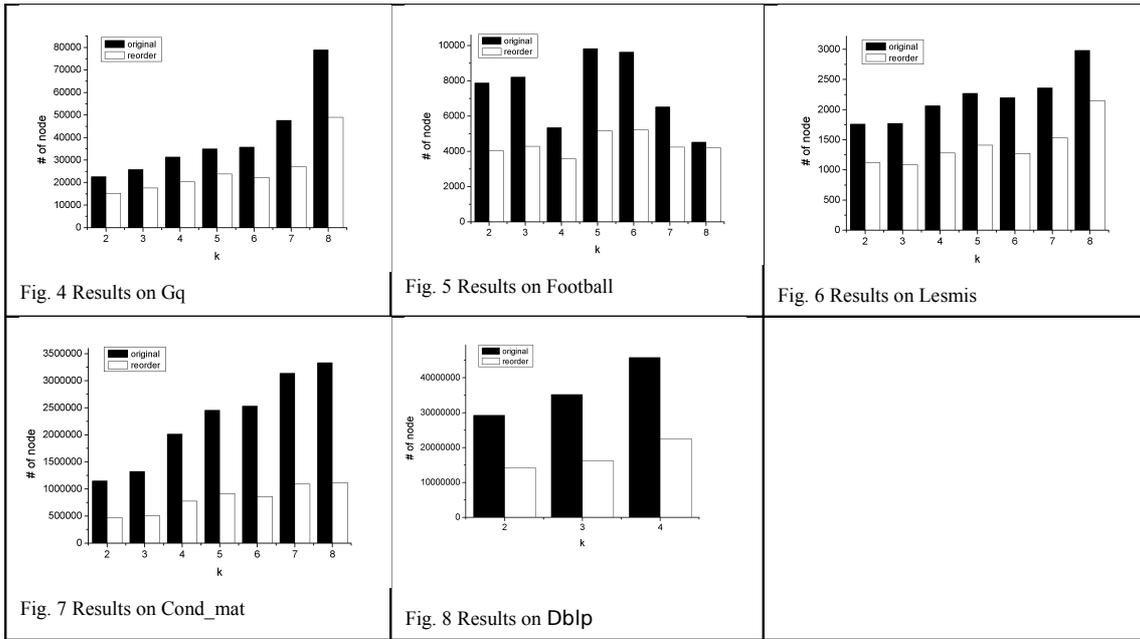
| Network         | $N$   | $M$    | $2M/N$ |
|-----------------|-------|--------|--------|
| <b>Gq</b>       | 578   | 1068   | 3.696  |
| <b>Lesmis</b>   | 77    | 254    | 6.597  |
| <b>Football</b> | 115   | 613    | 10.661 |
| <b>Cond_mat</b> | 36458 | 171736 | 9.421  |
| <b>Dblp</b>     | 48143 | 171932 | 7.143  |
|                 | 3     | 0      |        |

### 4.2. Results

Based on this specification, the following represents our experiment results.

#### 4.2.1 Re-sort of Vertices

This experiment is conducted to show the relative compression efficiency of a  $K^2$  tree constructed using the DFS heuristics, compared with that of a naïve  $K^2$  tree.



We compare the number of nodes in a  $K^2$  tree obtained by the original method and DFS, with a heuristic rule on the real network data listed on table 1, when  $K=2, 3, 4, 5, 6, 7$ . The experimental results are given in Figures 4 - 8. In Gq, Lesmis, Football, Cond\_mat and DBLP networks, the ratio of numbers of nodes of  $K^2$  tree using DFS, with heuristic rule and nodes of the original  $K^2$  tree is 64.22%, 63.40%, 62.14%, 36.73% and 47.8% on average. In those graphs, the number of nodes reduces to 42.07%, 38.0%, 48.76%, 66.5% and 54.0% respectively in the best case. These experiments indicate that the  $K^2$  tree using DFS and a heuristic rule can effectively reduce the number of nodes in a  $K^2$  tree.

Similarly, for different real networks, the optimal  $K$  value is also different. Generally, the number of nodes of a  $K^2$  tree will increase with the growth of  $K$ . However, there are some counter examples. For the Football network data, when  $K=4$ , the compression effectiveness reaches its optimum. It follows, therefore that the optimal  $K$  value depends on the structural property of a community networks.

#### 4.2.2 Adaptive $K$ compared with the fixed $K$

This experiment is conducted to illustrate the efficiency of an adaptive adjustment strategy for  $K$ , in comparison to a naïve  $K^2$  tree.

Table 2 Adaptive  $K$  compared with fixed  $K$

| Graph    | Original $K^2$ tree |        |        | $K^2$ tree with vertex ordered by DFS code |        |        |       |
|----------|---------------------|--------|--------|--|--------|--------|-------|
|          | $K=2$               | $K=3$  | $K=4$  | $K=2$                                      | $K=3$  | $K=4$  | SA    |
| Lesmis   | 1760                | 1771   | 2064   | 1120                                       | 1086   | 1280   | 452   |
|          | 25.68%              | 22.52% | 21.90% | 40.36%                                     | 41.62% | 35.31% |       |
| Football | 7876                | 8211   | 5340   | 4036                                       | 4280   | 3568   | 1471  |
|          | 18.68%              | 17.91% | 27.55% | 36.45%                                     | 34.37% | 41.23% |       |
| Gq       | 22598               | 25788  | 31365  | 15172                                      | 17587  | 20390  | 6667  |
|          | 29.50%              | 25.85% | 21.26% | 43.94%                                     | 37.91% | 32.70% |       |
| Cond_mat | 1151592             | 132314 | 201464 | 472716                                     | 508969 | 771004 | 17946 |
|          | 15.58%              | 9      | 0      | 37.96%                                     | 35.26% | 23.28% | 0     |
|          |                     | 13.56% | 8.91%  |  |        |        |       |

For a given node order, we compare the number of nodes of the self-adaptive  $k$  algorithm with that of the original  $K^2$  tree. In table 2, we present the experimental results for Lesmis, Football, Gq and Cond\_mat when  $K=2, 3, 4$ . Based on the DFS encoding with a heuristic rule, the number of nodes

in a self-adaptive  $K^2$  tree reduces to 39.10%, 37.35%, 38.18% and 32.17% of the original  $K^2$  tree on average. In the best case, the ratio is 35.31%, 34.37%, 32.70% and 23.28% respectively. So, by self-adaptively altering  $k$ , we can minimize damage to the community structure of the input graph. From Table 2 we can also observe that by combining the two optimization methods together, we can achieve a compression effectiveness of 21.90%, 17.91%, 21.26% and 8.91% in the best case.

### 4.2.3 Selection of $K$

This experiment is conducted to find the best  $K$  of an adaptive adjustment strategy when used to build  $K^2$  trees from real networks.

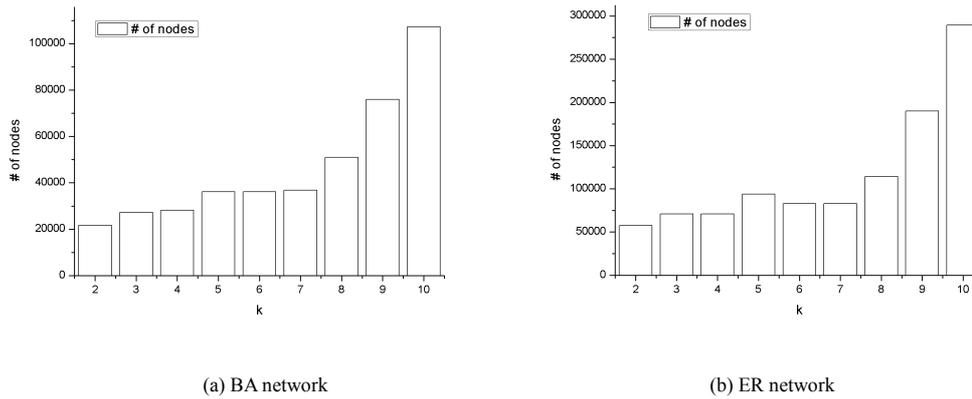


Fig. 9 Experimental result under different  $K$

An appropriate  $k$  value interval is of great importance to our compress algorithm. We use 2 synthetic networks to determine the best  $k$  value interval. The first BA network has 1000 nodes and 1980 edges and the second ER network has 1000 nodes and 5141 edges. For  $k$  varies from 2 to 10, the result on those two graph is given in figure 9. The number of nodes in  $K^2$  tree grows with the increment of  $k$  value. When  $k$  falls in  $[2,4]$ ,  $K^2$  tree has a acceptable result.

### 4.2.4 Influence of community structure

In this subsection, we will investigate how a community structure of a real network can impact upon on the efficiency of our approach.

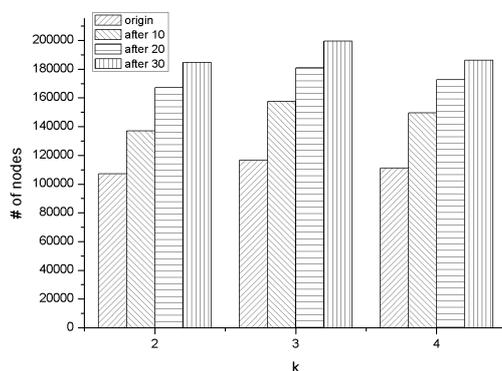


Fig. 10 Experimental result under different community structure

In this section we will discuss the impact of the clarity of a community structure in a network on the method we propose. We divide 1000 nodes into 10 communities. For a node pair  $(u,v)$ , if they

are in the same community, then the edge probability is 0.7, otherwise the probability is 0.001. As a result, we obtain a graph with 1000 nodes and 39124 edges. It is self-evident from the construction of synthetic networks that the resulting synthetic graphs are of strong community structure. In order to test the impact mentioned in this section, we do a  $t$  time experiment. In each experiment, we randomly disconnect  $10t$  edges inside a community and then randomly link  $10t$  edges between different communities. In our setting,  $t$  varies from 1 to 3. The result on those synthetic networks is given in Figure 10. It is shown that the compression grows worse with the loss of community structure. These results demonstrate that a mature NoP shall encompass two key characteristics. Firstly, the overall behavior of a mature NoP should be homogenous which in effect demonstrates that participant nodes interact with other well defined and trusted participant nodes. The results also demonstrate that such a behavior should be expected to occur during stage three of a NoP. This is because participant nodes have developed the skill, knowledge and expertise to process filtered requests. The results show also that the more explicit community structure is reflected in better compression. Secondly, the structure of the underlying network shall be broadly stable. That is to say that communication links between participant nodes will tend to be altered for the purposes of relationship maintenance, rather than as a general attempt to establish a relationship. Again, the results demonstrate that the more dynamic the structural changes, the worse the compression. It is expected that a newly formed CoP is more likely to exhibit a greater proportion of participant nodes joining and leaving (as per stage one) and thus will not tend to demonstrate stability.

## 5. Conclusion

A key thrust of this research has been to investigate how complex data graphs, such as those that represent NoP communities, can be appraised in terms of their relative maturity. As described earlier, we consider a network to be mature when two conditions are met; although individual nodes may possess heterogeneous behaviors, the overall behavior of the network would be regarded as homogenous. Secondly the network that represents a NoP will have progressed through three stages of maturity, based on those proposed by [23].

We have also described the simple  $K^2$  tree approach and illustrated how the natural occurrence of sparse regions in real-world networks can be exploited to achieve effective compression.

However, since this method does not consider the structural characteristics of a community network, there is still an opportunity to optimize compression further for graphs of NoP.

Therefore, we have proposed two optimization methods as follows: encoding nodes using DFS with heuristic rule, and a self-adaptive  $k$  building  $K^2$  tree. These two methods take good advantage of community feature in the real network, so as to reduce the number of internal nodes further.

Experiments on synthetic and real networks indicate that the  $K^2$  tree has less nodes than simple  $K^2$  tree. To conclude, we have chosen to exploit the observation that an augmented  $K^2$  tree compression method can be used to validate that a network has in fact reached stage three maturity.

## Acknowledgments

This work was supported by the National Natural Science Foundation of China under grants No.61003001; Specialized Research Fund for the Doctoral Program of Higher Education No.20100071120032; the Natural Science Foundation of Jiangsu Province No.BK2009153, BK2010280.

## References

- [1] comScore: <http://www.cctime.com/html/2010-4-22/201042291342486.htm>
- [2] Aggarwal C, Xie Y, Yu P S. Gconnect: A connectivity index for massive disk-resident graphs. Proc. of VLDB'09. Lyon, France: ACM, 2009: 862-873
- [3] Brisaboa N R, Ladra S, Navarro G.  $k^2$ -trees for compact web graph representation. LNCS 5721: Proc of the 16th Int Symp on String Processing and information Retrieval. Berlin: Springer, 2009: 18-30
- [4] Boldi P, Vigna S. The webgraph framework I: Compression techniques. Proc of the 13th Int Conf World Wide Web2004, New York, USA. ACM, 2004: 595-602
- [5] Boldi P, Vigna S. The Webgraph framework ii: Codes for the world-wide web. Proc of Data Compression Conference (DCC '04), Washington, USA: IEEE Computer Society, 2004: 528
- [6] Buehrer G, Chellapilla K. A scalable pattern mining approach to web graph compression with communities Proc of the 1st WSDM'08, California, USA: ACM, 2008: 95-106
- [7] Xiao Yanghua, MacArthur B D, Wang Hui, et al. Network Quotients: Structural Skeletons of Complex Systems. Physical Review E, 2008, 78: 046102
- [8] Chierichetti F, Kumar R, Lattanzi S, et al. On Compressing social networks. Proc of the SIGKDD'09. Paris, France, 2009: 219-228
- [9] Li Yuming, Dong Wanpeng, Peng Yinghong. Study on Matrix Compressive Storage Method Based on 0-1 Property-matrix. Computer Engineering and Applications. 2003,39(2):82-84 (in Chinese)
- [10] Sergio L. Toral, M. Rocio Martinez-Torres, Federico Barrero and Francisco Cortes. An empirical study of the driving forces behind online communities. Internet Research, 2009, 19(4): 378-392
- [11] M.R. Martinez-Torres, S.L. Toral, F. Barrero and F. Cortes. The role of Internet in the development of future software projects. Internet Research, 2010, 20(1): 72-86
- [12] E. Wenger, Communities of Practice: Learning, Meaning, and Identity, Cambridge University Press, Cambridge, 1998
- [13] J.P. Johnson, Collaboration, peer review and open source software, Information Economics and Policy 18 (2006) 477-497.
- [14] Samet H. The Quadtree and Related Hierarchical Data Structures. ACM Computing Surveys, 1984, 16 (2): 187-260
- [15] Garey M R, Johnson D S, Stockmeyer L. Some simplified NP-complete graph problems Proc the 6th Annual ACM Symp Theoret Comput. Sci.,1974,1(3):237-267,197
- [16] Faloutsos M, Faloutsos P, and Faloutsos C. On power-law relationships of the internet topology. Proc. of SIGCOMM. Cambridge MA, USA. ACM, 1999: 251-262
- [17] Newman M E J. The structure and function of complex networks. SIAM Review, 2003, 45(2):167-256
- [18] Broder A Z, Glassman S C, Manasse M S, et al. Syntactic clustering of the Web. Computer Networks, 1997, 29(8-13):1157-1166
- [19] Broder A, Charikar M, Frieze A, Mitzenmacher M. Min-wise independent permutations. Journal of Computer and System Sciences, 2000, 60:630-659
- [20] Gibson D, Kumar R, Tomkins A. Discovering Large Dense Subgraphs in Massive Graphs. Proc. of VLDB'05. Trondheim, Norway: ACM, 2005. 721-732
- [21] Knuth D E. The Stanford GraphBase: A Platform for Combinatorial Computing. [ACM](http://www.acm.org) Press, 1993
- [22] Batagelj V, Mrvar A. (2006): Pajek datasets. <http://vlado.fmf.uni-lj.si/pub/networks/data/>
- [23] Newman M E J. The structure of scientific collaboration networks. Proc Natl Acad Sci USA. 2001, 98(2): 404-409
- [24] Dblp [2010-05-5] <http://dblp.uni-trier.de/xml/>