

## 用多层次聚类法完成的大规模关系图的可视化\*

黄茂林<sup>+</sup>, NGUYEN Quang Vinh

(Faculty of Information Technology, University of Technology, Sydney, Australia)

### Large Graph Visualization by Hierarchical Clustering

HUANG Mao-Lin<sup>+</sup>, NGUYEN Quang Vinh

(Faculty of Information Technology, University of Technology, Sydney, Australia)

+ Corresponding author: E-mail: maolin@it.uts.edu.au

**Huang ML, Nguyen QV. Large graph visualization by hierarchical clustering. *Journal of Software*, 2008,19(8): 1933–1946. <http://www.jos.org.cn/1000-9825/19/1933.htm>**

**Abstract:** This paper proposes a new technique for visualizing large graphs of several ten thousands of vertices and edges. To achieve a graph abstraction, a hierarchical clustered graph is extracted from a general large graph based on the community structures discovered in the graph. An enclosure geometrical partitioning algorithm is then applied to achieving the space optimization. For graph drawing, it uses a combination of spring-embedder and circular drawing algorithms that archives the goal of optimization of display space and aesthetical niceness. The paper also discusses an interaction mechanism accompanied with the layout solution. The interaction not only allows users to navigate hierarchically through the entire clustered graph, but also provides a way to navigate multiple clusters concurrently. Animation is also implemented to preserve user mental maps during the interaction.

**Key words:** graph drawing; information visualization; view navigation; interaction; clustered graph

**摘要:** 提出了一种新的大规模图形可视化技术.它可显示含有几万个接点和边的大规模关系图.为了完成对图形的抽象化,一个多层次的聚类图形从原始的大规模关系图中抽取了出来.这种抽取是建立在大规模关系图的内在结构基础上来完成的.一种递归封装式的几何划分算法被应用来完成对几何空间的优化,在具体的制图技术上,使用了一种用力导向布局算法和环形制图法相结合的新方法,从而完成了对显示空间的优化和美学上的优化.同时也讨论了相关的人机交互技术,所采用的人机交互算法不仅能让使用者从上到下层次式地浏览整个聚类图形,同时也能提供多层次聚类图形的并行浏览.动画技术也同时被运用,以保护使用者的精神图不被打乱.

**关键词:** 图形绘制;信息可视化;场景游览;人机交互;聚类图形

中图分类号: TP391      文献标识码: A

## 1 Introduction

Graph visualization has been widely used in human-computer interaction. A graph commonly includes a node set and an edge set to represent entities and relationships between entities respectively. Graphs generated in

real-world applications could be very large with thousands or perhaps millions of nodes, such as citation and collaboration networks and the World Wide Web (WWW). As the result of rapid increasing of the size in networks, the large scale visualization has become one of the hottest topics in Information Visualization. The question about how to comprehensively display large graphs on the screen becomes the key issue in graph visualization. However, the display of large graphs can decrease significantly the performance of a visualization technique which normally performs well on small or medium size of datasets. Large graph visualization usually suffers from poor running time and the limitation of display space. In addition, the issue of "view-ability" and usability also arises because it will be almost impossible to discern between nodes and edges when a dataset of thousands of items are displayed<sup>[1]</sup>.

It seems that classical graph models with a simple node-link diagram tend to be inadequate for large scale visualization with several thousands of items. The lack of formal hierarchical structures in real world applications could limit the conveying and perception of the complicated information. Figure 1 shows an example of the graph visualization of a WWW site which illustrates two typical major problems:

- Too many nodes (pages) to be displayed and the layout of such a large geometrical area could not be fitted in one single screen
- The layout of the graph has inefficient utilization of display space with many unused areas in the display.

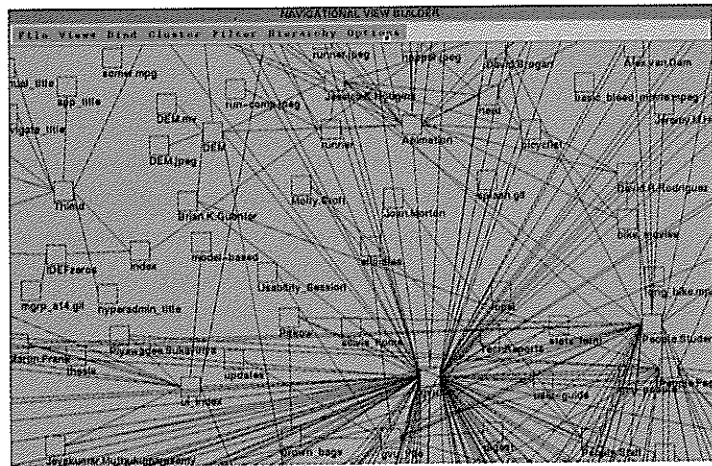


Fig.1 An example of a large graph visualization using the classic virtual-page technique

To address the first problem, a well established new graph model to accommodate with the visualization of large graphs is required. We believe that one way to deal with the display of large graphs is to provide users with a certain degree of Graph Abstract. That is to filter out some details of the graph drawing which is assumed at a time the viewer is not interested, while the overall structure of the graph drawing is maintained for navigation.

Among several available graph visualization approaches, we believe that the use of clustered graph is a better option for graph abstraction. Therefore, a good visualization system for very large graphs should be a combination of three components including graph drawing, graph clustering and interaction<sup>[2,3]</sup>. Visualization of clustered graphs such as the ones in Refs.[4,5] is one of an excellent approaches to deal with large graphs through the graph abstraction. A clustered graph can be extracted from a general graphs by partitioning recursively the graph into a hierarchy of sub-graphs, so that it simplifies the complex structure of the large graph for easy interpretation, perception and navigation of large information spaces.

To solve the second problem mentioned above, we need to optimize layout algorithms to maximize the utilization of display screen by allowing more nodes to be displayed. The research from Ware<sup>[6]</sup> shows that more information can be displayed on very high-resolution and large screen, but it does not necessarily provide very

much more information into the brain. This is because the conventional monitor covers only 5–10% of visual field in the normal condition, but it uses as much as 50% of brain pixels<sup>[7]</sup>. The study also shows that the uniquely stimulated brain pixels peak at the width of a normal monitor view, and it is effective (but not critical) to increase the number of pixels for the normal desktop to reach the limit of the brain pixels. Therefore, investigation of optimized visual abstraction (clustering) techniques that could provide viewers with more comprehensive views of the large graphs becomes important.

## 2 Related Work

Large graphs visualization has recently received a lot of attention from researchers in both information visualization and graph drawing communities. Although some newly available techniques such as techniques found at<sup>[3,8–12]</sup>, are quite capable of visualizing large graphs of thousands to hundred thousands of nodes and edges, visualization of large graphs, is still one of the open topics in information visualization.

Harel and Koren<sup>[8]</sup> described a technique to draw a graph that used high-dimensional embedding and then projected it onto a 2D plane. Although this technique is very fast and is capable of exhibiting graph in various dimensions with some good navigational ability, it is more suitable for visualizing mesh-graphs rather than tree-like graphs or clustered graphs. One of the good approaches for handling large graphs is to use multi-scale visualization<sup>[11,13,14]</sup>. This approach typically applies a force-directed algorithm to draw large graphs using multi-scale scheme in which they try to beautify the coarsest-scale representation. Techniques in this approach aim to improve the processing speed while maintaining the graph niceness. A good visualization of large graphs can also be achieved by using multilevel techniques<sup>[9,15]</sup>. In short, these techniques improve the visual appearance of the visualization by defining different levels for a structure so that they can present the graphs using an optimal algorithm at each level.

Although the above techniques are quite capable of visualizing large graphs, the space-efficiency is not considered in the visualization which could limit the amount of information to be visualized on the screen at a time. Fekete, *et al.*<sup>[16]</sup> presented a space-efficient visualization of graph using a modification of the well-known Tree-Maps<sup>[17]</sup>. Technically, the authors used Tree-Maps to display the tree structure of graph and used explicit link curves to present the other links. This technique was optimal in term of using display space and it is quite useful for visualizing structures that the underlying trees have some meaning. However, it did not perform well in general graphs and clustered graphs because the link curves might cause unnatural look of the graphs. Some preliminary works have been carried out and from which two tree visualization techniques *Space-Optimized Tree* (or *SO-Tree*)<sup>[20]</sup> and *EncCon Tree*<sup>[19]</sup> have been developed that can quickly display large trees with maximized utilization of display space. However, these solutions are only suitable for trees (hierarchical structures).

This paper proposes a new technique for visualizing very large general graphs. Our proposed technique is very similar to the framework of Tulip<sup>[3]</sup> which consists of three components: *graph clustering*, *graph layout* and *interaction*. We first use a new clustering algorithm to partition the complete graph into abstract clusters for achieving the view abstraction; that aims to reduce the visual complexity of the graph layout, and enhance the comprehension and understanding of the graph.

The clustered graph is then visualized using a new space-efficient layout technique which is a combination of different layout algorithms. This geometrical optimization of graph layout allows more data items (and clusters) to be displayed within a limited screen resolution. Our visualization provides viewers with not only an abstract view of the entire graph but also an interaction technique for the navigation of large graphs. The navigation method allows users to browse hierarchically through the clustered graph and navigate across a number of selected clusters. All the

interactions are accommodated with animation to preserve user mental maps during the navigation.

### 3 The Architecture of the Visualization

Our model for visualizing large graphs includes several processes which are illustrated at Fig.2. There are two major phases involved in this model including the *clustering analysis* and the *user interface* phases. The two phases operate independently.

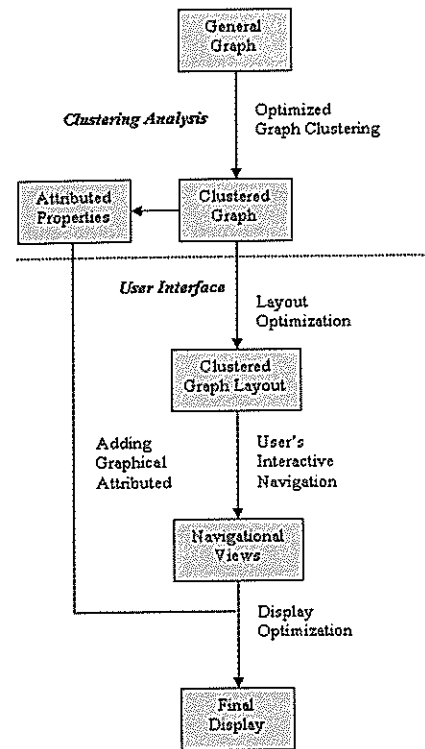


Fig.2 A architectural model for visualizing large graphs

layout and navigation algorithms is required for handling hundred thousands of items within minutes or seconds using a personal computer with limited display space and computational power.

The final display is created through the view navigation and graphical properties. We use rich graphic attributes to assist viewers to quickly identify the domain specific properties associated with data items. We next describe briefly of the clustering and the technical detail of our visualization technique.

### 4 Graph Clustering

We use a graph clustering method which can quickly discover the community structure embedded in a large graph and divide the graph into densely connected sub-graphs. The graph clustering algorithm partitions the graphs into smaller sub-graphs based on the density of connection within and between subgroups. Although this process does not require a very fast algorithm, the computational cost of clustering algorithms should be controlled with the worst-case running time  $O(n^2)$  or faster on a sparse graph to ensure its capability of handling hundreds thousands of items within a few hours using an ordinary personal computer.

The *clustering analysis* phase is responsible for analyzing a large graph and partitioning it into a clustered graph based on discovered internal communities. In short, the clustering algorithm recursively divides the graphs into smaller sub-graphs based on the density of connection within and between subgroups. Although this process does not require a very fast algorithm and it can operate independently, the computational cost of clustering algorithms should be controlled with running time of  $O(n^2)$  or better on a sparse graph to ensure its capability of handling hundreds thousands of items within a few hours using an ordinary personal computer. The clustering process also extracts attributed properties for nodes, edges and relations between sub-graphs.

The *user interface* phase is responsible for the visualization and navigation of the clustered graph, including layout optimization, interactive viewing and display optimization. A combination of different layout algorithms is employed which aims to optimize the geometrical space and so that the large graph can be drawn at a normal screen size.

During the navigation of clustered graphs, we allow users to interactively adjust the views to reach an optimal representation of the graph; from where they can obtain the best understanding of the data and structures. This visualization is involved with real time human-computer interaction. Therefore, very fast graph

Research in graph clustering for large datasets has recently received a lot of effort from researchers. However, the discovery of fast and effective clustering algorithms for handling large graphs is still a big challenge. In fact, the finding of an exact solution for graph clustering is still believed to be an NP-complete problem. Kernighan and Lin<sup>[23]</sup> and Newman and Girvan<sup>[24]</sup> have presented their heuristic techniques that can produce quite good solutions for graph clustering. However, their techniques are very slow with the worst-case running time  $O(m^2n)$  or  $O(n^3)$  on sparse graphs. This makes it almost impossible to partition a graph with more than a few thousands of elements. Newman<sup>[18]</sup> later proposed a new fast algorithm for detecting community structure in networks. This method runs in worst-case time  $O((m+n)n)$  or  $O(n^2)$  on a sparse graph. The algorithm is very fast and can handle graphs with hundred thousands of elements. However, its clustering results are not very consistent, especially a poor balance between clusters.

Although the quality or usefulness of an embedded graph drawing algorithm is highly dependent on its application domain, aesthetics is still one of the most important quality factors in graph drawing or graph visualization in which the readability of a graph is measured or justified. The aesthetical criteria for graph drawing can be found from a book on Graph Drawing by Di Battista, *et al.*<sup>[25]</sup> and the revised version from Ware, *et al.*<sup>[26]</sup>. Among the aesthetical criteria, the even distribution of vertices and the maximization of display symmetry of a graph structure are two important factors to ensure the quality of graph visualization. These criteria are closely related to the quality of the balance of clusters produced by a clustering algorithm. In other words, we believe that a good clustering algorithm should achieve the goals of balanced clustering in which in each level of the hierarchy the size of the clusters should be about the same. This property helps associated graph drawing method to provide good visualization in term of readability.

This paper used a graph clustering method<sup>[27]</sup> which can quickly discover the community structure embedded in large graphs and divide the graph into densely connected sub-graphs. The proposed algorithm can not only run fast in time  $O((m+n)n)$ , but also achieve a consistent partitioning result in which a graph is divided into a set of clusters of the similar size. Although our objective, i.e. keeping clusters balance, is similar to Duncan, *et al.*<sup>[28]</sup>, this clustering algorithm is more general rather than just by using a binary space partition (BSP) clustering.

The balance in size of clusters provides users with a clearer view of the clustered graph and thus it makes it easier to visualize and navigate large graphs. Our balanced clustering technique creates a layout optimization at both global and local levels of the display through the use of enclosure+connection layout technique. This allows more visual items to be displayed within limited screen resolutions and with comprehensive views. The combination between our clustering method and a space-efficient layout technique would enable the visualization of very large general graphs with several thousands of elements. Figure 3 shows an example of a clustering output using the balanced clustering algorithm on a large and highly connected graph.

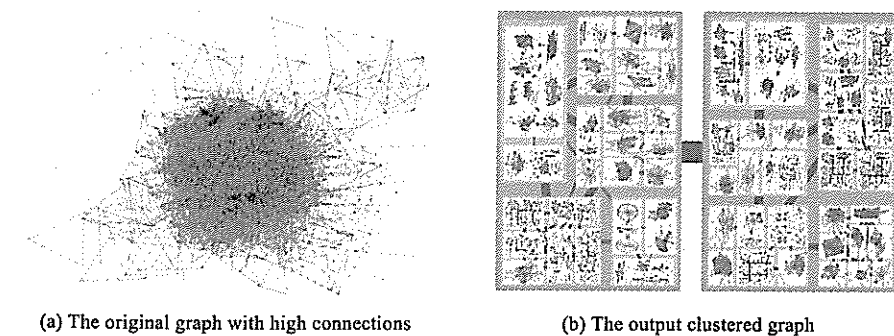


Fig.3 An example of a clustering

## 5 Graph Visualization

We use a new space-efficient visualization technique similarly to *EncCon*<sup>[19]</sup> to optimize the geometrical space for visualizing large clustered graphs with several thousands of nodes and edges. This technique consists of two components, the space-efficient layout and the interactive navigation. The layout of clustered graph is generated by using a combination of an extended fast enclosure partitioning algorithm, called *Clenccon*, and a number of traditional graph drawing algorithms, including a *spring-embedder* algorithm<sup>[21]</sup>, a circular drawing, and simple layout algorithm, to archive the objectives of space-efficiency, aesthetical niceness and fast computation. Although some existing techniques can use any layout algorithm at each clusters, such as using a *spring-force* algorithm<sup>[9,11]</sup>, in our belief, the choice of a space-efficient enclosure-partitioning algorithm at high levels will be more efficient because it provides more space for displaying information at the limited display.

The *Clenccon* layout algorithm is only applied to those non-leaf sub-graphs in which the space utilization and computational cost are crucial. In other cases, the other layout algorithms, including *spring-embedder* and *circular drawing*, is applied to the calculation of the position for those leaf sub-graphs, which contain a small number of nodes in which the space utilization issue becomes less important and, therefore, the aesthetic niceness and flexibility issues need to be further considered. The use of a particular layout algorithm depends on the nature of the leaf sub-graphs. Our system also displays a high-level node-link diagram to present the overall clustering structure explicitly (see examples at Figs. 5 and 6).

Our *Clenccon* layout algorithm inherits essentially the advantage of space-filling techniques<sup>[17,19]</sup> that utilize display space by using area division for the partitioning of sub-trees and nodes. Note that the issue of space utilization becomes significantly important when visualizing large graphs with thousands or even hundred thousands of nodes and edges because of the limitation of screen pixels. It is similar to *EncCon*<sup>[19]</sup> that uses a rectangular division method for recursively positioning the nodes hierarchically. This property aims to provide users with a more straightforward way to perceive the visualization and ensures the efficient use of display space. However, our new technique is applied for clustered graphs rather than simple tree structures. Therefore, the algorithm takes the connectivity property between sibling nodes into its partitioning process. We now describe the technical detail of our layout and navigation algorithms.

### 5.1 Layout algorithm

Our layout algorithm is responsible for positioning of all nodes in a given clustered graph in a two-dimensional geometrical space, including a vertex subset  $\{v_1, v_2, \dots, v_n\}$  in  $V$  and a cluster subset  $\{v'_1, v'_2, \dots, v'_n\}$  in  $V'$ . A clustered graph  $C = \{G, T\}$  is derived by a general graph  $G = \{V, E\}$  and a cluster tree  $T$  whose leaves are in  $V$ . Each cluster  $v'_i = C$  is a sub clustered graph, contains a subset of  $V$  given by the leaves of the sub-tree  $T'$  rooted at  $v'$ . The root  $v'$  of the sub-tree  $T'$  is also called a super-node. The super nodes are not displayed in our visualization but they are used for partitioning process of calculating the local region for sub clustered graph. For the partitioning of clustered graph  $C$ , We define a virtual tree consisting of a set of super-nodes for area division. We define a super-node  $r(v'_i)$  for each cluster  $v'_i$ . Further description of the clustered graph can be found at<sup>[5]</sup>.

The layout algorithm is a combination of two algorithms: 1) *Clenccon* - a fast area division algorithm and 2) graph drawing algorithms, including a *spring-embedder* algorithm, a *circular drawing* and a simple algorithm to lay out a very small number of nodes. Each cluster  $v'_i$  is bounded by a rectangular local region  $R(v'_i)$  centered at super-node  $r(v'_i)$  and the drawing of the corresponding sub-clustered-graph  $G(v'_i)$  is restricted to be inside the geometrical area of  $R(v'_i)$ . Therefore, the local region  $R(v'_i)$  of cluster  $v'_i$  is the sum of the rectangular areas assigned to its children. The position of the super-node  $r(v'_i)$  of  $v'_i$  is at the centre of the rectangle defined by  $R(v'_i)$ . The

position of leaf nodes is defined by either the *spring embedder*, the *circular drawing* or the simple layout algorithms.

We first assign the entire rectangular display area as the local region to the clustered graph  $C$ . We then recursively partition the local regions for every sub-clusters until all the clusters are reached.

We assign a weight  $w(v')$  to each vertex  $v'$  for the calculation of the local region  $R(v')$  of the vertex. Although the weight of each vertex can be associated with its property, all the leaf vertices in our experiments have the same weight. Suppose that the rectangular local region  $R(v')$  for cluster  $v'$  is drawn, we then need to calculate the local regions  $\{R(v'_{i+1}), R(v'_{i+2}), \dots, R(v'_{i+k})\}$  for its sub-clusters  $\{v'_{i+1}, v'_{i+2}, \dots, v'_{i+k}\}$ . The partitioning ensures that the area of each rectangle  $R(v'_{i+1})$  is proportional to the weight  $w(v'_{i+1})$  of the cluster  $v'_{i+1}$ . The calculation of  $w(v')$  of a cluster  $v'$  is done recursively from leaves of the cluster tree to the root of the cluster tree. The calculation is done by the following formula:

$$w(v) = w_o + S \sum_{i=1}^k w(v_{i+1}),$$

where  $w_o$  is the internal weight of cluster  $v'$ . Although the internal weight of a cluster can be defined by the cluster's attributed property, we define the internal weight of all clusters to be 1 for all experiments.  $S$  is a constant ( $0 < S < 1$ ), and  $w(v'_{i+1})$  is the weight assigned to the  $i$ th child of cluster  $v'$ . The constant  $S$  determines the size difference of local regions of all clusters based on the number of descendants of those vertices.

The process of recursive partitioning  $R(v')$  into sub-regions  $\{R(v'_{i+1}), R(v'_{i+2}), \dots, R(v'_{i+k})\}$  for all its children clusters  $\{v'_{i+1}, v'_{i+2}, \dots, v'_{i+k}\}$  is illustrated as the procedure below:

#### procedure partitioning (Node $N$ )

```
{
  if all child-nodes of  $N$  are leaf-nodes then
  {
    lay out the child-nodes using a graph algorithm;
    scale the layout to fit with rectangular local region;
  }
  else
  {
    lay out the child-nodes using Clenccon algorithm;
    for each non-leaf child-node of  $N$ 
    {
      partitioning(child-node);
    }
  }
}
```

#### procedure Clenccon-layout (Node [] Nodes)

```
{
  group linked-nodes into subgroups;
  sort the subgroups based on connection and size;
  lay out subgroups using EncCon algorithm;
```

```

for each subgroup
{
    lay out nodes in subgroup using EncCon algorithm;
}
}

procedure childnode-layout (Node [] Nodes)
{
    if number of child-nodes < K1 then
    {
        lay out the child-nodes using simple algorithm;
    }
    else if number of child-nodes > K1 and
        no. edges / no. child-nodes > K2 then
    {
        lay out the child-nodes using circular drawing;
    }
    else
    {
        lay out the child-nodes using Spring algorithm;
    }
}

```

where  $K1=6$  indicates the number of nodes that is suitable for each algorithm. If there are just a few number of nodes, i.e. less than 10 nodes, a simple node location algorithm can perform well.  $K2=5$  indicates the ratio of the number of edges over the number of nodes. Because the force directed algorithms do not usually perform well for a graph whose the number of edges is much larger than number of nodes, we use the circular drawing in this situation.

The detail description of the area partitioning *EncCon* can be found at<sup>[19]</sup>. Although there are several improved force-directed layout algorithm, the traditional *Spring Embedder*<sup>[21]</sup> layout algorithm is chosen in our implementation. This is because the algorithm is simple, easy to implement, flexible and perform well in general for a small number of nodes in which a more complicated algorithm is not necessary. The circular drawing is a simple technique but is very effective to show the pattern of relationship for a graph whose the number of edges is much larger than number of nodes. Technically, we place all nodes equally on a circle where relational nodes are located close together so that the pattern of connections can be display more clearly (see Figs.5 and 6).

Figure 4(a) shows an example of partitioning and drawing a small clustered graph using our algorithm. We can see that the algorithm uses the *Clenccon* algorithm to layout three cluster nodes and their inter-relationships to ensure the efficient utilization of space and uses the *spring-embedder* algorithm to draw the sub-graphs within each cluster to achieve the aesthetic niceness and flexibility. Note that the inter-relationships among clusters here are represented by using abstract links. Figure 4(b) shows the same example of the partitioning, but the inter-relationships among clusters are represented by using the original structure of relationships. Figure 5 to Figure 7 illustrate the visualizations of our layout algorithm on various very large datasets. These pictures show clearly the structure of clustered graphs, in which sub-graphs are efficiently partitioned and drawn inside their local regions. All of these pictures use abstract links to represent the inter-relational structures among clusters.

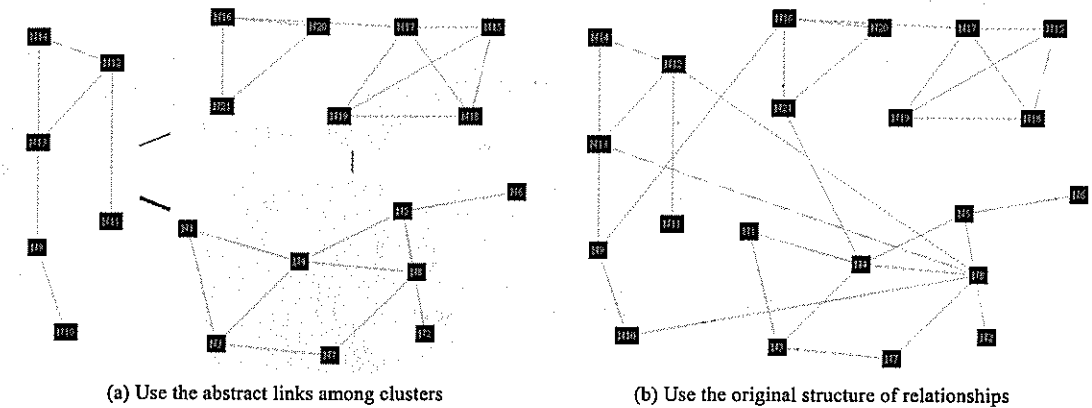


Fig.4 An example of partitioning and drawing a small clustered graph using our algorithm

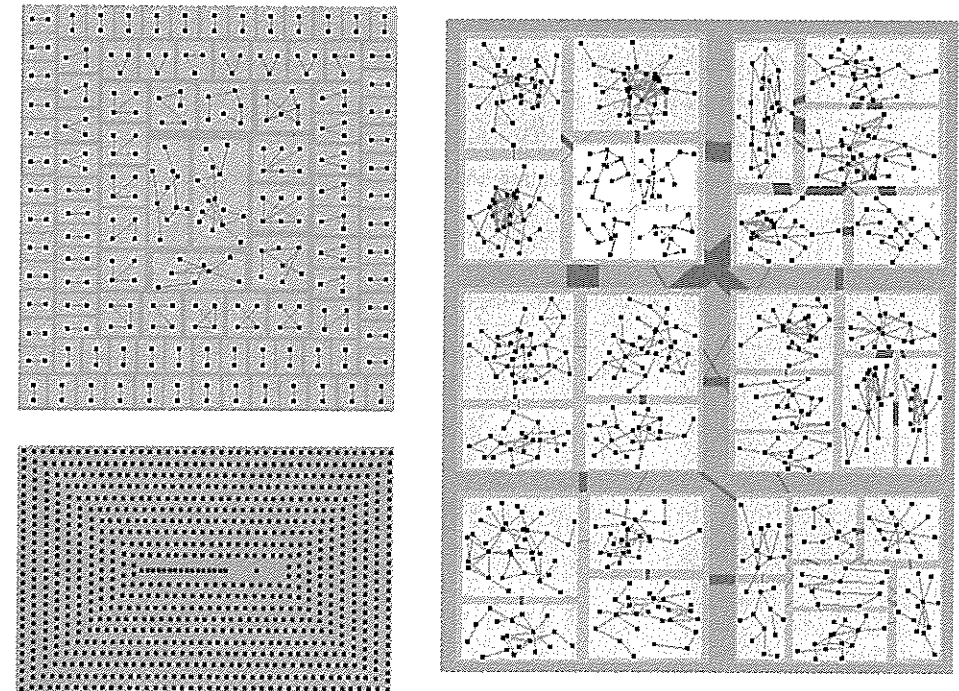


Fig.5 An example of clustering on a large citation dataset with over 2 000 vertices and 4 200 edges of the social network papers

## 5.2 Navigational views

There is no pure visualization technique that could assist data retrieval without providing users with an associate navigation mechanism in graphic user interface design. In our user interface component, during the navigation we enable users to interactively adjust the views to reach an optimized representation of the graph; from which users can obtain the best understanding of the data and its relational structures they are interesting in.

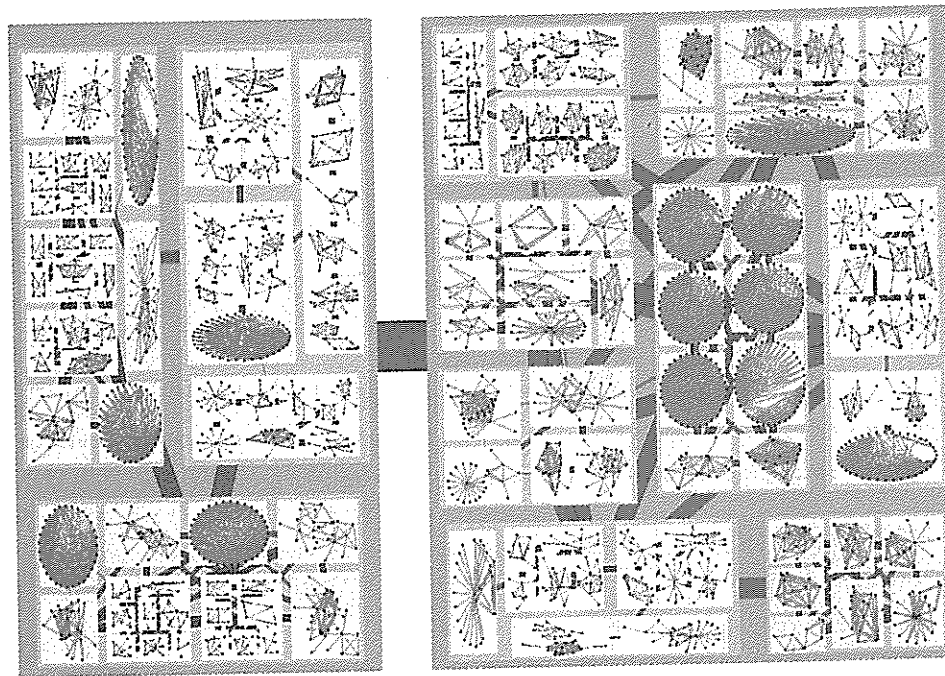


Fig.6 An example of clustering on a protein dataset with very strong relations produced by our algorithm

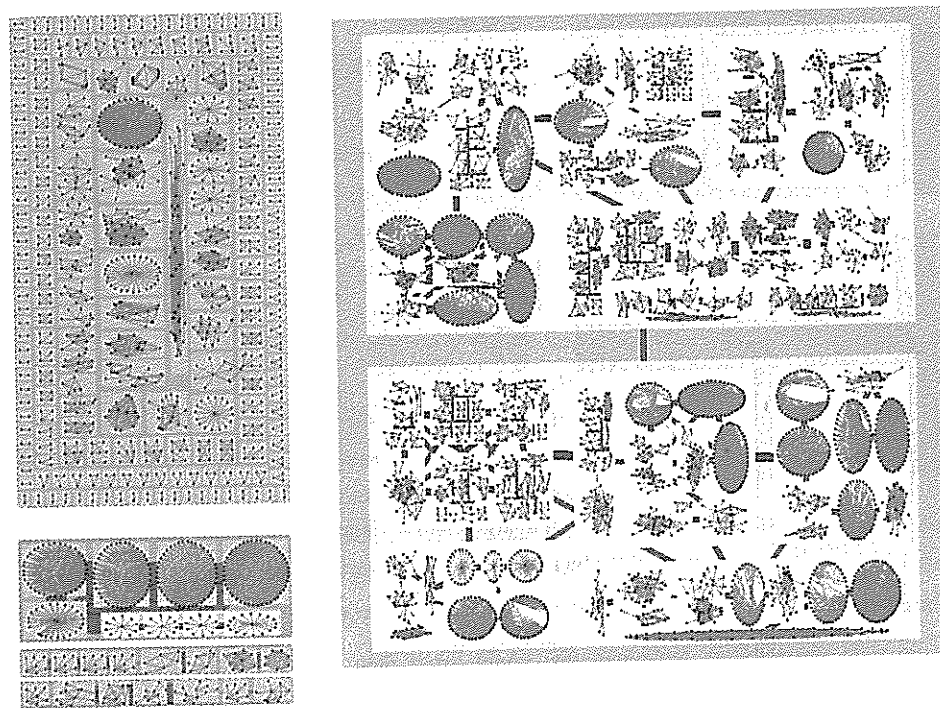


Fig.7 An example of clustering on a protein dataset with our algorithm cluster a protein network with over 15 000 nodes and 40 000 edges

In our prototype, we use a *multiple-views* technique<sup>[22]</sup>, including a *main view* and *context views*, to achieve the *focus+context* view navigation of large clustered graphs. The *main view* shows as much detail of information as

possible allowing users to efficiently perform their interaction and visual analysis on this area. The *context views*, which are displayed in the small areas at the left-hand side, are only responsible for displaying a several levels of the contextual information (or the history of navigation) during the navigation. Therefore, a large amount of the detail in these context views is filtered and they remain only the main structures and important landmarks in an abstract manner for guiding the navigation. These simplified structures could preserve user's mental map of where they are, where they come from, and where they have been during the navigation.

The navigational views allow the exploration of data hierarchically as well as across multiple selected clusters to quickly focus on the interest parts of data. Therefore, users can semantically zoom-in one or many areas of interest or sub-clusters while retaining simplified context views. It effectively uses both focus and context views to enable the interactive investigation of very large data. We provide three views to assist viewers to obtain the best understanding of the dataset and its relational structure they are currently exploring. These views are: 1) a *full context view*, 2) a *current context views* and 3) a *main view* (see Figs.8 and 9).

*Full context view* is displayed as a small panel located at the top-left side of the visualization. This view displays the entire context of a large clustered graph in high-level abstraction with little details. This enables users to always maintain an overall view of the information. The *full-context view* only shows three levels of the clustered graph from the root in a simplified display. This is because too much information displaying at a small area would create the overcrowded and distraction to users. The view updates automatically according to changes occurred in the database. It highlights the context of the current sub-graph the user is interacting, so that users can always identify the where they are and where they have been in a large information space. For example, in Figure 8 we can easily identify the *current-context view* which is on the top-left region of the *full-context view* and the *main-view* which is the right sub-graph inside the top-left region of the *full-context view*.

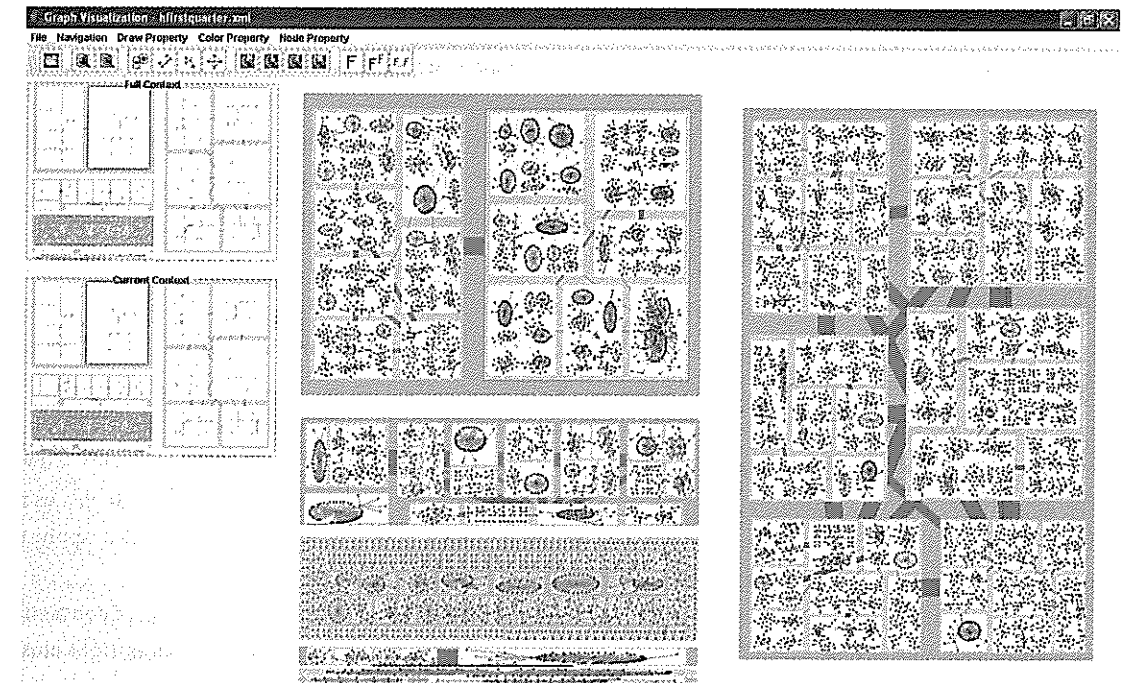


Fig.8 A clustered visualization of a very large protein dataset with three views: 1) a full-context, 2) a current-context view and 3) a main view

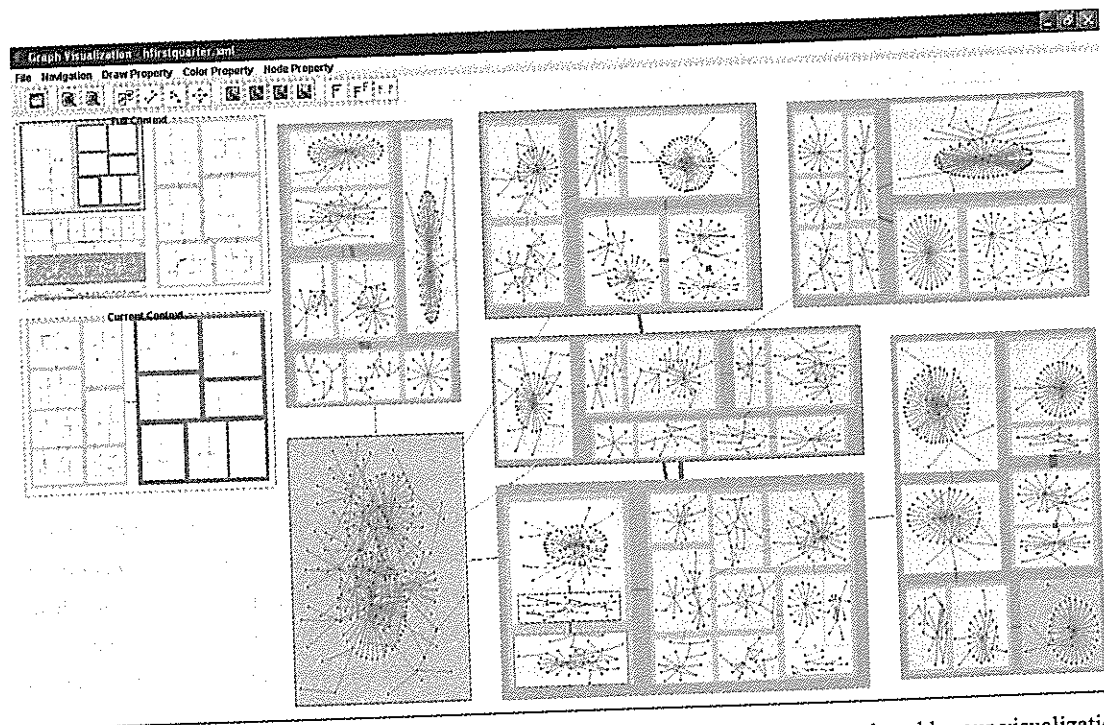


Fig.9 A navigational view of a selected sub-graph of the same data set as Fig.8 produced by our visualization

*Current context view* is displayed as a small panel locating under the *full-context view* at the left-hand side. This view displays the immediate context or the up-level view of the main view. Similar to the *full-context view*, this view displays up to three levels of clustered graph from the root in an abstraction manner. The *current-context* panel also updates automatically and highlights a focused sub-graph which to be viewed in the main view in details. This enables users to easily identify their focused region at the current context (see Fig.9).

*Main view* is displayed in a large area of the visualization. It displays the detail of a focused sub-graph. It displays the structure of the clustered graph in an enclosure manner with color brightness. This view also shows the connectivity strength between clusters using the width of edges. Furthermore, the size of nodes in the main-view is automatically rescale based on the number of visible nodes.

### 5.3 Final display

Some graphical attributes are employed to in final display of clustered graphs to assist viewers to quickly identify the domain specific properties of data and the hierarchical structure of the clustered graph. Colors are employed to assist viewers to quickly identify the hierarchical structure of the clustered graph. In our prototype, the local regions of nodes at different levels are painted with a same color but at different brightness. This drawing property aims to provide a pleased view while retaining the clarity between sub-graphs. Although the use of brightness of colors for background can theoretically apply to several levels of hierarchies, we also apply this drawing property to first four hierarchical levels to avoid the confusion with too many colors.

Width or thickness of the edge is employed to represents the weight of the edge (or the number of connections between two nodes). For example, in Fig.4(b), there are 3 links between the left cluster and bottom-right cluster. However, Fig.4(a) displays only one thick abstract link between these clusters. We can see that edges among leaf nodes are drawn with light-green color and edges among non-leaf nodes (or between two clusters) are drawn with light-gray color. Furthermore, to avoid the overlapping due to the over-thickness of edges, the thickness of an edge

is limited. Therefore, if the weight of an edge between two clusters is greater than a limited number in our implementation, the edge is drawn at its limited maximum-width with darker boundary.

Although theoretically the area partitioning algorithm can optimize display space of the entire graph layout, a small portion of display space is reserved as gaps among clusters for showing abstract links between clusters and these abstract links indicate the connectivity strength between clusters.

## 6 Conclusions

We have presented a new visualization technique that could handle the visualization of very large graphs with up to tens thousands or even hundred thousands of elements on an ordinary Personal Computer. Our method includes two independent steps: clustering and visualization. The clustering step aims to reduce the visual complexity and enhance the comprehension of large graph layouts through the use of visual abstraction. We first discover an optimized community structure in a graph and divide it into densely connected clusters. We then use three levels of visual abstraction: 1) the full context view, 2) the current context view and 3) the main view, to display the large graph. The visualization uses a new space-efficient layout and navigation technique that can visualize effectively the large clustered graphs with several thousands of elements on a limited display space. We use a *multiple view* technique to archive the *focus+context* view navigation. The interaction Animation is also employed to preserve the user mental maps during the interaction.

## References:

- [1] Herman I, Melancon G, Marshall MS. Graph visualization in information visualization: A survey. *IEEE Trans. on Visualization and Computer Graphics*, 2000,(6):24-44.
- [2] Marshall S. Methods and tools for the visualization and navigation of graphs [Ph.D. Thesis]. Bordeaux: University Bordeaux I, 2001.
- [3] Auber D. Tulip: a huge graph visualization framework. In: Mutzel P, Junger M, eds. *Graph Drawing Software, Mathematics and Visualization*. Springer-Verlag, 2003, 105-126.
- [4] Eades P, Feng Q. Multilevel visualization of clustered graphs. In: North SC, ed. *Graph Drawing (GD'96)*. California: Springer, 1996. 101-112.
- [5] Feng Q. Algorithms for drawing clustered graphs [Ph.D. Thesis]. Newcastle: University of Newcastle, 1997.
- [6] Ware C. *Information Visualization: Perception for Design*. San Francisco: Morgan Kaufmann Publishers, 2004.
- [7] Wilkins AJ. *Visual Stress*. Oxford: Oxford University Press, 1995.
- [8] Harel D, Koren Y. Graph drawing by high-dimensional embedding. *Journal of Graph Algorithms and Applications*, 2004,8(2): 195-214.
- [9] Walshaw C. A multilevel algorithm for force-directed graph drawing. *Mathematics Research Report 00/IM/60*, University of Greenwich, 2000.
- [10] Abello J, van Ham F, Krishnan N. ASK-GraphView: A large scale graph visualization system. *IEEE Trans. on Visualization and Computer Graphics*, 2006,12(5):669-676.
- [11] Harel D, Koren Y. A fast multi-scale method for drawing large graphs. In: Marks J, ed. *Graph Drawing*. Colonial Williamsburg: Springer-Verlag, 2000. 183-196.
- [12] Hachul S, Junger M. An experimental comparison of fast algorithms for drawing general large graphs. In: Healy P, Nikolov NS, ed. *Graph Drawing*. Limerick: Springer-Verlag, 2005. 235-250.
- [13] Auber D, Chiricota Y, Jourdan F, Malancon G. Multiscale visualization of small world networks. In: *IEEE Symp. on Information Visualization*. Seattle: IEEE Computer Society, 2003. 75-81.
- [14] Auber D, Jourdan F. Interactive refinement of multi-scale network clusterings. In: *Int'l Conf. on Information Visualization*. London: IEEE, 2005. 703-709.

- [15] Archambault D, Munzner T, Auber D. TopoLayout: Multi-Level graph layout by topological features. *IEEE Trans. on Visualization and Computer Graphics*, 2007,13(2):305-317.
- [16] Fekete JD, Wang D, Dang N, Aris A, Plaisant C. Overlaying graph links on treemaps. In: *IEEE Symp. on Information Visualization - Symp. Poster Compendium*. Seattle: IEEE Computer Society, 2003. 82-83.
- [17] Johnson B, Shneiderman B. Tree-Maps: A space-filling approach to the visualization of hierarchical information structures. In: Kaufman AE, ed. *1991 IEEE Visualization*. IEEE Computer Society, 1991. 284-291.
- [18] Newman MEJ. Fast algorithm for detecting community structure in networks. *Physical Review E*, 2004,69:066133.
- [19] Nguyen QV, Huang ML. EncCon: An approach to constructing interactive visualization of large hierarchical data. *Information Visualization Journal*, 2005,4(1):1-21.
- [20] Nguyen QV, Huang ML. Space-Optimized tree: A connection+enclosure approach for the visualization of large hierarchies. *Information Visualization Journal*, 2003,2(1):3-15.
- [21] Eades P. A heuristic for graph drawing. *Congressus Numerantium*, 1984,42:149-160.
- [22] Robert JC. On encouraging multiple views for visualisation. In: *Int'l Conf. on Information Visualization*. London: IEEE, 1998. 8-14.
- [23] Kernighan G, Lin S. An efficient heuristic procedure for partitioning graphs. *The Bell System Technical Journal*, 1970,29(2): 291-307.
- [24] Newman MEJ, Girvan M. Finding and evaluating community structure in networks. *Physical Review E*, 2004,69:026113.
- [25] di Battista G, Eades P, Tamassia R, Tollis IG. *Graph Drawing: Algorithms for the Visualization of Graphs*. Prentice Hall, 1999.
- [26] Ware C, Purchase H, Colpoys L, McGill M. Cognitive measurements of graph aesthetics. *Information Visualization Journal*, 2002, 1(2):103-110.
- [27] Huang ML, Nguyen QV. A fast algorithm for balanced graph clustering. In: *Proc. of the Int'l Conf. on Information Visualization*. Zürich: IEEE, 2007. 46-52.
- [28] Duncan CA, Goodrich MT, Kobourov SG. Balanced aspect ratio trees and their use for drawing large graphs. *Journal of Graph Algorithms and Applications*, 2000,4(3):19-46.



**HUANG Mao-Lin** is a senior lecturer and doctoral supervisor at the Faculty of Information Technology, University of Technology, Sydney, Australia. His research areas are graph drawing and information visualization.



**NGUYEN Quang Vinh** is a Post-Doctoral Fellow at the Faculty of Information Technology, University of Technology, Sydney, Australia. His current research interests include information visualization.





中国科协精品科技期刊工程  
项目资助期刊



中国计算机学会会刊



科学出版社

# 软件学报

(Ruanjian Xuebao)

第 19 卷第 8 期

2008 年 8 月

## 目次

### 软件可视化和信息可视化

- 软件可视化和信息可视化专刊前言 ..... 张康 吕建 (1865)
- 一种面向图形化建模语言表示法的元模型 ..... 何啸 麻志毅 邵维忠 (1867)
- 基于图文法的动态软件体系结构支撑环境 ..... 马晓星 曹春 余萍 周宇 (1881)
- 一种基于边的上下文相关图文法形式化框架 ..... 曾晓勤 韩秀清 邹阳 (1893)
- 可视化语言技术在软件开发中的应用(英文) ..... 孔骏 赵春颖 (1902)
- 一种用于 Marching-Graph 图形绘制的快速收敛布局算法(英文) ..... 全武 黄茂林 (1920)
- 用多层次聚类法完成的大规模关系图的可视化(英文) ..... 黄茂林 NGUYEN Quang Vinh (1933)
- 一种模型驱动的交互式信息可视化开发方法 .....  
..... 任磊 王威信 周明骏 滕东兴 马翠霞 戴国忠 王宏安 (1947)
- 数据聚类中基于浓度噪音消除的可视化参数选择方法(英文) ..... 钱宇 (1965)
- 基于社会网络可视化分析的数据挖掘(英文) ..... 杨育彬 李宁 张瑶 (1980)
- 三维人体运动特征可视化与交互式运动分割 ..... 肖俊 庄越挺 吴飞 (1995)
- 地质断层三维可视化模型的构建方法与实现技术 ..... 朱良峰 潘信 吴信才 刘修国 (2004)

### 数据库技术

- 数据空间技术研究 ..... 李玉坤 孟小峰 张相於 (2018)
- XML 数据流上的高效聚集算法 ..... 王宏志 李建中 骆吉洲 (2032)
- 语义查询扩展中词语-概念相关度的计算 ..... 田莹 杜小勇 李海华 (2043)
- 基于查询采样的高维数据混合索引 ..... 张军旗 周向东 施伯乐 (2054)
- 数据流处理中确定性 QoS 的保证方法 ..... 武珊珊 于戈 吕雁飞 谷峪 李晓静 (2066)
- XML 流数据查询结果的缓存管理 ..... 杨卫东 王清明 施伯乐 (2080)

### 计算机网络与信息安全

- 传感器网络中基于虚拟坐标的节点调度方案 ..... 李小龙 林亚平 易叶青 余建平 卢新国 (2089)
- 基于规则推导的特权隐式授权分析 ..... 蔡嘉勇 卿斯汉 刘伟 何建波 (2102)
- 一种基于关键属性的优化数据一致性维护方法 ..... 周婧 王意洁 李思昆 (2114)
- 无线局域网 EDCA 机制 MAC 接入延时分析 ..... 周新运 皇甫伟 孙利民 王显雷 李凯慧 (2127)
- 网络入侵检测中的自动决定聚类数算法 ..... 肖立中 邵志清 马汉华 王秀英 刘刚 (2140)
- 基于混合跳链条件随机场的异构 Web 记录集成方法 ..... 黄健斌 姬红兵 孙鹤立 (2149)
- 基于 Weil 对的多接收者公钥加密方案 ..... 鲁力 胡磊 (2159)

关于推荐 2008 年 CCF 优秀博士学位论文的通知 ..... (2158)

《软件学报》投稿指南 ..... (封三)

# 软件学报 *Software*

---

## Journal of *Software*

---



中国科学院软件研究所 主办  
中国计算机学会  
科学出版社 出版

Sponsored by Institute of Software, The Chinese Academy of Sciences  
and the China Computer Federation  
Published by Science Press

Beijing China

## Editor-in-Chief

LI Ming-Shu(李明树) *Institute of Software, The Chinese Academy of Sciences (CAS), China*

## Associate Editors-in-Chief

LU Ru-Qian(陆汝钤) *Institute of Mathematics, AMSS, The Chinese Academy of Sciences (CAS), China*

Yao, Andrew C. *Tsinghua University, China*

## Editorial Board

Bjørner, Dines, *Technical University of Denmark, Denmark*

CHEN Huo-Wang(陈火旺), *National University of Defense Technology, China*

CHEN Ke-Fei(陈克非), *Shanghai Jiaotong University, China*

Chin, Francis Yuk-Lun, *University of Hong Kong, China*

Clarke, Edmund M., *Carnegie Mellon University, USA*

DAI Guo-Zhong(戴国忠), *Institute of Software, CAS, China*

DONG Yun-Mei(董璠美), *Institute of Software, CAS, China*

FENG Deng-Guo(冯登国), *Institute of Software, CAS, China*

FENG Yu-Lin(冯玉琳), *Institute of Software, CAS, China*

Graham, Ronald L., *University of California at San Diego, USA*

HE Xin-Gui(何新贵), *Peking University, China*

HUANG Tao(黄涛), *Institute of Software, CAS, China*

JIN Zhi(金芝), *Institute of Mathematics, AMSS, CAS, China*

JU Jiu-Bin(鞠九滨), *Jilin University, China*

Krieg-Bruckner, Bernd, *Universitat Bremen, Germany*

LI De-Yi(李德毅), *Institute of Electronic System Engineering, China*

LI Guo-Jie(李国杰), *Institute of Computing Technology, CAS, China*

LI Hua(李华), *Institute of Computing Technology, CAS, China*

LI Jian-Zhong(李建中), *Harbin Institute of Technology, China*

Li, Kai, *Princeton University, USA*

LI Lei(李磊), *Sun Yat-Sen University, China*

Li, Ming, *University of Waterloo, Canada*

LI Wei(李未), *Beijing University of Aeronautics and Astronautics, China*

LI Xuan-Dong(李宣东), *Nanjing University, China*

LIN Chuang(林闯), *Tsinghua University, China*

LIN Hui-Min(林惠民), *Institute of Software, CAS, China*

LIU Chun-Nian(刘椿年), *Beijing University of Industry, China*

LIU Ji-Ren(刘积仁), *Northeastern University, China*

LU Wei-Ming(陆维明), *Institute of Mathematics, AMSS, CAS, China*

LÜ Jian(吕建), *Nanjing University, China*

MEI Hong(梅宏), *Peking University, China*

Motiwalla, Juzar, *National University of Singapore, Singapore*

PENG Qun-Sheng(彭群生), *Zhejiang University, China*

QIAN Hua-Lin(钱华林), *Computer Network Information Center, CAS, China*

SHI Bai-Le(施伯乐), *Fudan University, China*

SHI Chun-Yi(石纯一), *Tsinghua University, China*

SUN Jia-Guang(孙家广), *Tsinghua University, China*

SUN Zhong-Xiu(孙钟秀), *Nanjing University, China*

TAN Tie-Niu(谭铁牛), *Institute of Automation, CAS, China*

TIAN Jie(田捷), *Institute of Automation, CAS, China*

WANG Ji(王戟), *National University of Defense Technology, China*

WANG Shan(王珊), *People's University of China, China*

WANG Xiao-Yun(王小云), *Tsinghua University, China*

WU Jian-Ping(吴建平), *Tsinghua University, China*

XU Zhuo-Qun(许卓群), *Peking University, China*

YANG Fu-Qing(杨芙清), *Peking University, China*

ZHANG Bo(张钹), *Tsinghua University, China*

ZHAO Qin-Ping(赵沁平), *Beijing University of Aeronautics and Astronautics, China*

ZHENG Nan-Ning(郑南宁), *Xi'an Jiaotong University, China*

ZHENG Wei-Min(郑伟民), *Tsinghua University, China*

ZHOU Chao-Chen(周巢尘), *Institute of Software, CAS, China*

ZHOU Long-Xiang(周龙骧), *Institute of Mathematics, AMSS, CAS, China*

ZHOU Zhi-Hua(周志华), *Nanjing University, China*

Zhu, Hong, *Oxford Brookes University, UK*

SOFTWARE V

1865

1867

1881

1893

1902

1920

1933

1947

1965

1980

1995

2004

DATABASE T

2018

2032

2043

2054

2066

2080

COMPUTER M

2089

2102

2114

2127

2140

2149

2159

## Contents

**SOFTWARE VISUALIZATION AND INFORMATION VISUALIZATION**

- 1865 Preface of the Special Issue on Software Visualization and Information Visualization  
*ZHANG Kang, LÜ Jian*
- 1867 A Metamodel for the Notation of Graphical Modeling Languages  
*HE Xiao, MA Zhi-Yi, SHAO Wei-Zhong*
- 1881 A Supporting Environment Based on Graph Grammar for Dynamic Software Architectures  
*MA Xiao-Xing, CAO Chun, YU Ping, ZHOU Yu*
- 1893 An Edge-Based Context-Sensitive Graph Grammar Formalism  
*ZENG Xiao-Qin, HAN Xiu-Qing, ZOU Yang*
- 1902 Visual Language Techniques for Software Development  
*KONG Jun, ZHAO Chun-Ying*
- 1920 Fast Convergence Layout Algorithm for Drawing Graphs in Marching-Graph  
*QUAN Wu, HUANG Mao-Lin*
- 1933 Large Graph Visualization by Hierarchical Clustering  
*HUANG Mao-Lin, NGUYEN Quang Vinh*
- 1947 A Model Driven Development Method for Interactive Information Visualization  
*REN Lei, WANG Wei-Xin, ZHOU Ming-Jun, TENG Dong-Xing, MA Cui-Xia, DAI Guo-Zhong, WANG Hong-An*
- 1965 A Visual Approach to Parameter Selection of Density-Based Noise Removal for Effective Data Clustering  
*QIAN Yu*
- 1980 Networked Data Mining Based on Social Network Visualizations  
*YANG Yu-Bin, LI Ning, ZHANG Yao*
- 1995 Feature Visualization and Interactive Segmentation of 3D Human Motion  
*XIAO Jun, ZHUANG Yue-Ting, WU Fei*
- 2004 Construction Method and Actualizing Techniques of 3D Visual Model for Geological Faults  
*ZHU Liang-Feng, PAN Xin, WU Xin-Cai, LIU Xiu-Guo*

**DATABASE TECHNOLOGY**

- 2018 Research on Dataspace  
*LI Yu-Kun, MENG Xiao-Feng, ZHANG Xiang-Yu*
- 2032 Efficient Aggregation Algorithms on XML Stream  
*WANG Hong-Zhi, LI Jian-Zhong, LUO Ji-Zhou*
- 2043 Computing Term-Concept Association in Semantic-Based Query Expansion  
*TIAN Xuan, DU Xiao-Yong, LI Hai-Hua*
- 2054 High Dimensional Hybrid Index Based on Query Sampling  
*ZHANG Jun-Qi, ZHOU Xiang-Dong, SHI Bai-Le*
- 2066 A Deterministic QoS Guaranteeing Approach for Data Stream Processing  
*WU Shan-Shan, YU Ge, LÜ Yan-Fei, GU Yu, LI Xiao-Jing*
- 2080 Buffer Management of the Results of Queries over XML Streams  
*YANG Wei-Dong, WANG Qing-Ming, SHI Bai-Le*

**COMPUTER NETWORKS AND INFORMATION SECURITY**

- 2089 A Node Scheduling Scheme Based on Virtual Coordinate in Sensor Networks  
*LI Xiao-Long, LIN Ya-Ping, YI Ye-Qing, YU Jian-Ping, LU Xin-Guo*
- 2102 Analysis on Implicit Authorization in Privilege Through Rule Deduction  
*CAI Jia-Yong, QING Si-Han, LIU Wei, HE Jian-Bo*
- 2114 An Optimistic Data Consistency Maintenance Method Based on Key-Attributes  
*ZHOU Jing, WANG Yi-Jie, LI Si-Kun*
- 2127 MAC Access Delay Analysis of EDCA Mechanism of Wireless LANs  
*ZHOU Xin-Yun, HUANGFU Wei, SUN Li-Min, WANG Xian-Lei, LI Kai-Hui*
- 2140 An Algorithm for Automatic Clustering Number Determination in Networks Intrusion Detection  
*XIAO Li-Zhong, SHAO Zhi-Qing, MA Han-Hua, WANG Xiu-Ying, LIU Gang*
- 2149 Integration of Heterogeneous Web Records Using Mixed Skip-Chain Conditional Random Fields  
*HUANG Jian-Bin, JI Hong-Bing, SUN He-Li*
- 2159 Multi-Recipient Public Key Encryption Scheme Based on Weil Pairing  
*LU Li, HU Lei*

## 《软件学报》专刊征文通知

专刊名称: 医学影像处理与分析

特约编辑: 田捷(中科院自动化所)、白净(清华大学)、包尚联(北京大学)

出版日期: 2009年第2期

<http://www.jos.org.cn/news/MIPA.htm>

征文通知及投稿方式已发布在我刊网站上 <http://www.jos.org.cn>, 欢迎作者踊跃投稿。

软件学报

Ruanjian Xuebao

(月刊, 1990年创刊)

第19卷 第8期 2008年8月

Journal of Software

(monthly)

(Started in 1990)

Vol.19 No.8 Aug. 2008

编辑 《软件学报》编辑部  
(北京 8718 信箱 邮编 100190)  
电话: 010-62562563  
E-mail: [jos@iscas.ac.cn](mailto:jos@iscas.ac.cn)  
<http://www.jos.org.cn>

主编 李明树

主办单位 中国科学院软件研究所 中国计算机学会

主管单位 中国科学院

出版 科学出版社  
(北京东黄城根北街16号 邮编 100717)

印刷 北京北林印刷厂

总发行处 北京报刊发行局

订购处 全国各地邮局

国外总发行 中国国际图书贸易总公司  
(北京 399 信箱 邮编 100044)

Edited by Editorial Board of Journal of Software  
(P.O.Box 8718, Beijing 100190, P.R.China)  
E-mail: [jos@iscas.ac.cn](mailto:jos@iscas.ac.cn)  
<http://www.jos.org.cn>

Editor-in-Chief: LI Ming-Shu

Sponsored by Institute of Software, The Chinese Academy of Sciences and China Computer Federation

Published by Science Press (16 Donghuangchenggen North Street, Beijing 100717, China)

Printed by Beijing Bei Lin Printing House

Generally Distributed by Beijing Bureau for Distribution of Newspapers and Journals

Domestically Distributed by All Local Post Offices in China

Overseas Distributed by China International Book Trading Corporation (P.O.Box 399, Beijing 100044, China)

ISSN 1000-9825

CN 11-2560/TP

国内邮发代号: 82-367

国外发行代号: M4628

©2008 Journal of Software(版权所有)

定价: 36.00 元

公开发 行