

COMP5048

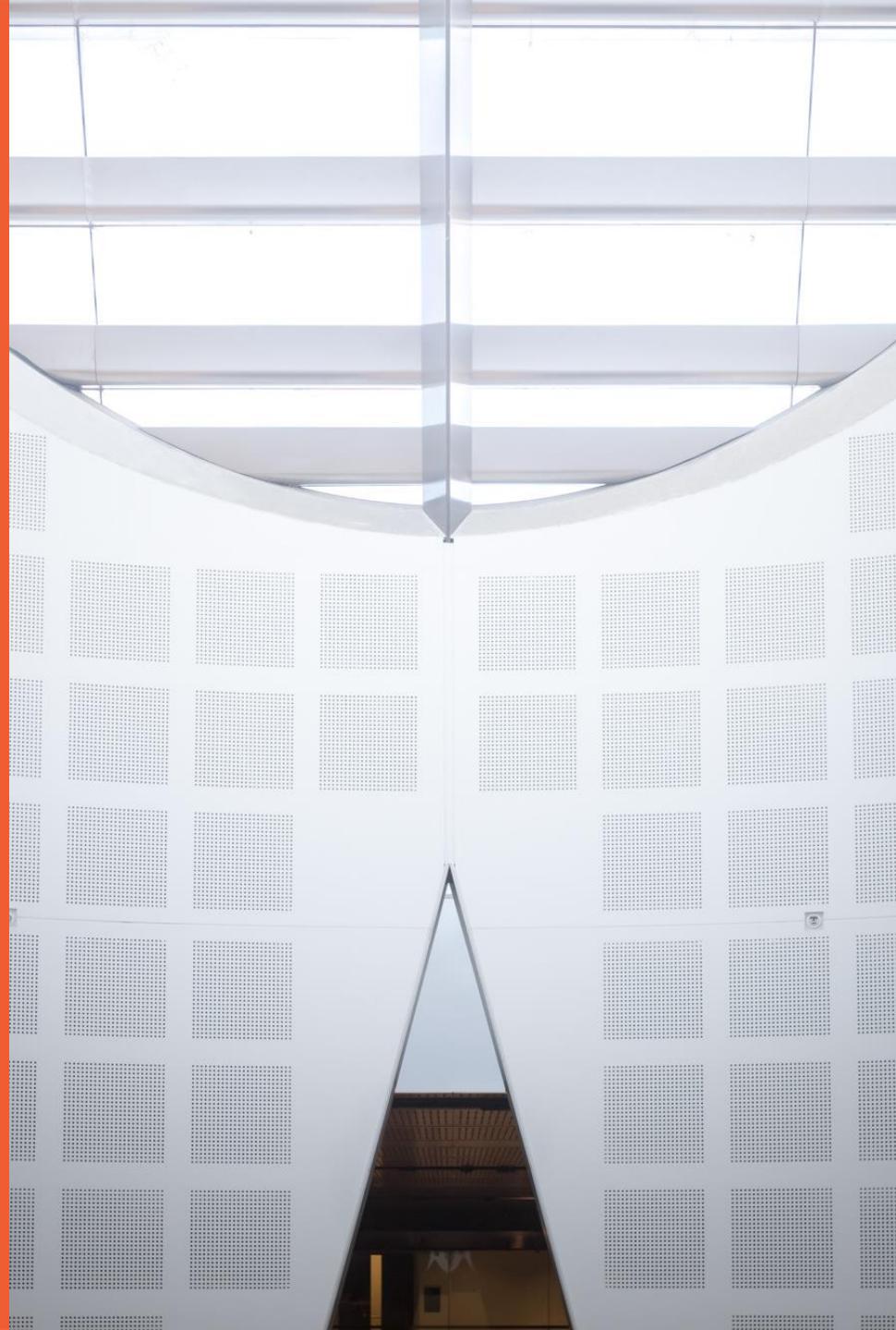
Visual Analytics

Week 5: Big Data Visual Analytics

Professor Seokhee Hong
School of Information Technologies



THE UNIVERSITY OF
SYDNEY



Copyright warning

COMMONWEALTH OF AUSTRALIA

Copyright Regulations 1969

WARNING

This material has been reproduced and communicated to you by or on behalf of the University of Sydney pursuant to Part VB of the Copyright Act 1968 (**the Act**).

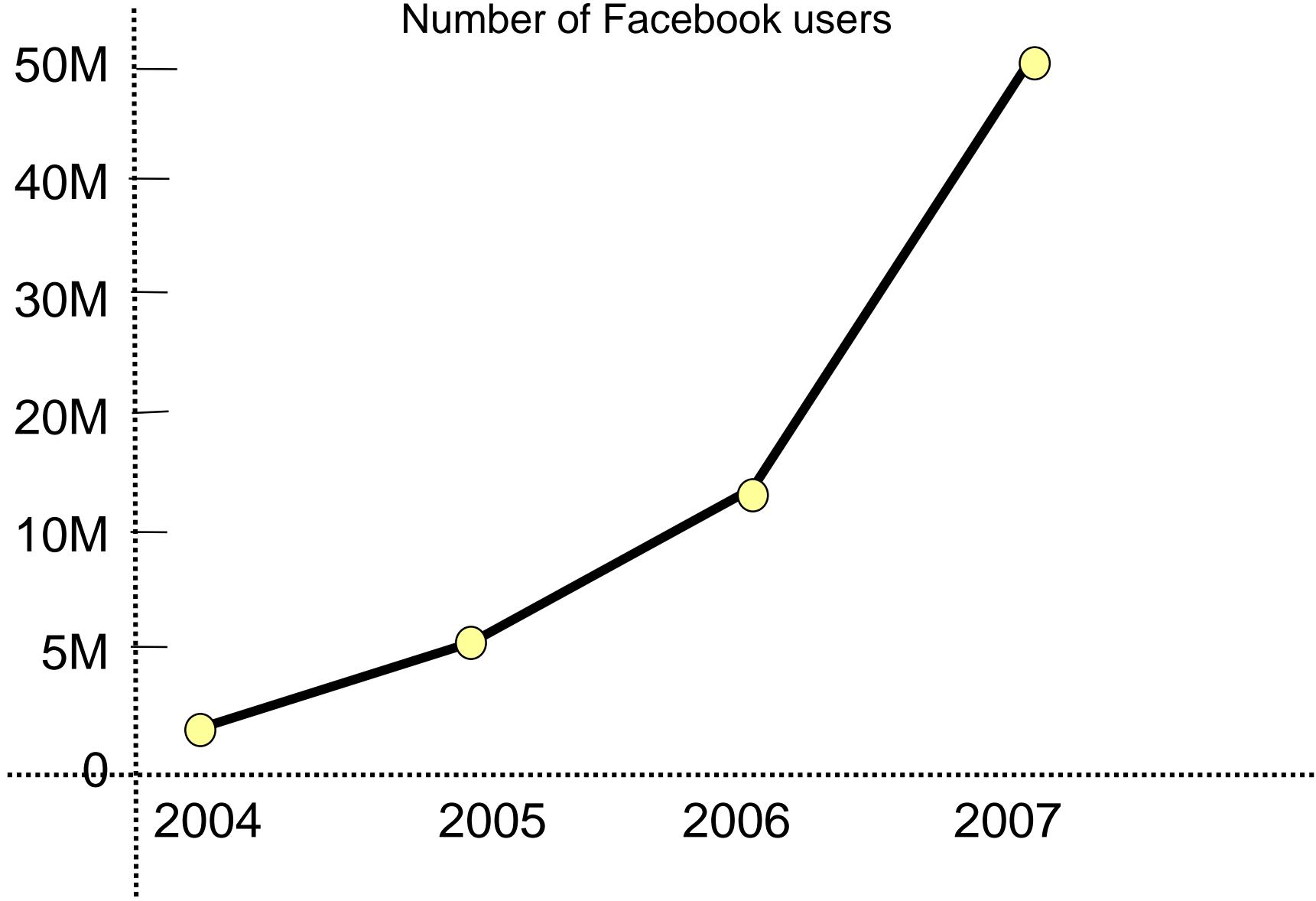
The material in this communication may be subject to copyright under the Act. Any further copying or communication of this material by you may be the subject of copyright protection under the Act.

Do not remove this notice.

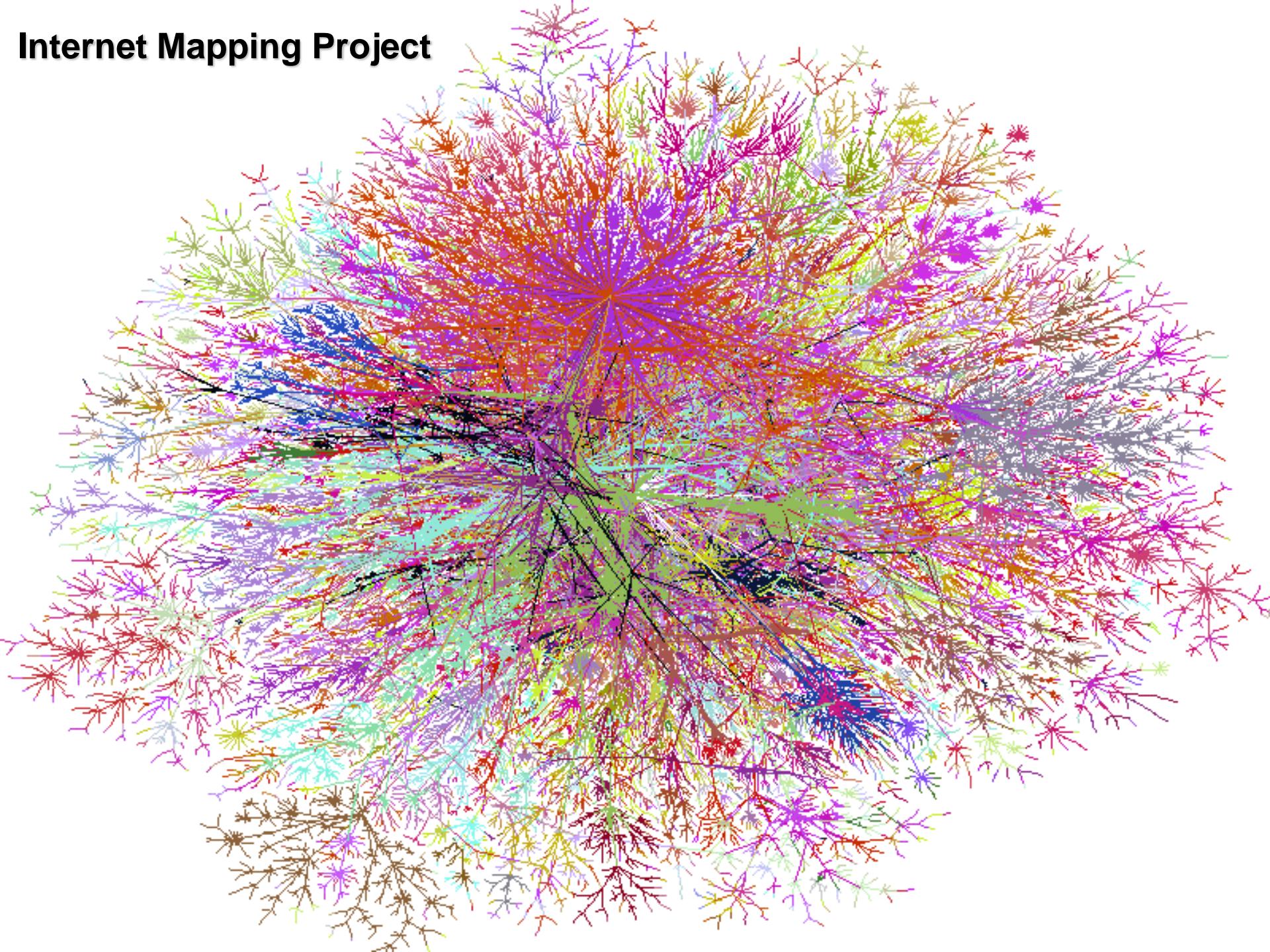
The Scale Problem

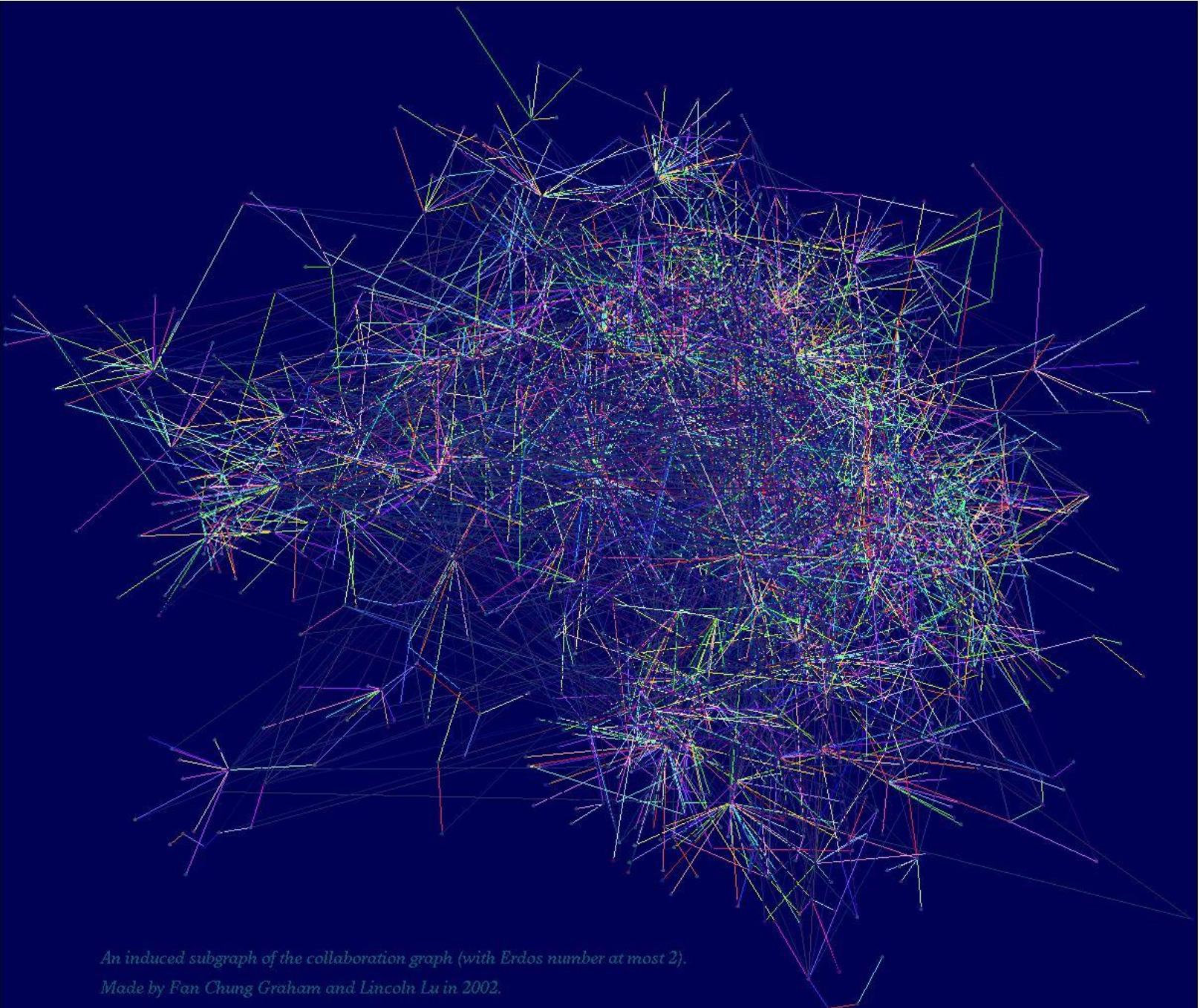
- Big Complex Data Challenge
- Data sets are growing at a faster rate than the human ability to understand them.
- Existing analysis/visualization methods do not scale well enough to be effective on current data sets.
- Human has limitation in perception and cognition.

Social networks:
Number of Facebook users



Internet Mapping Project



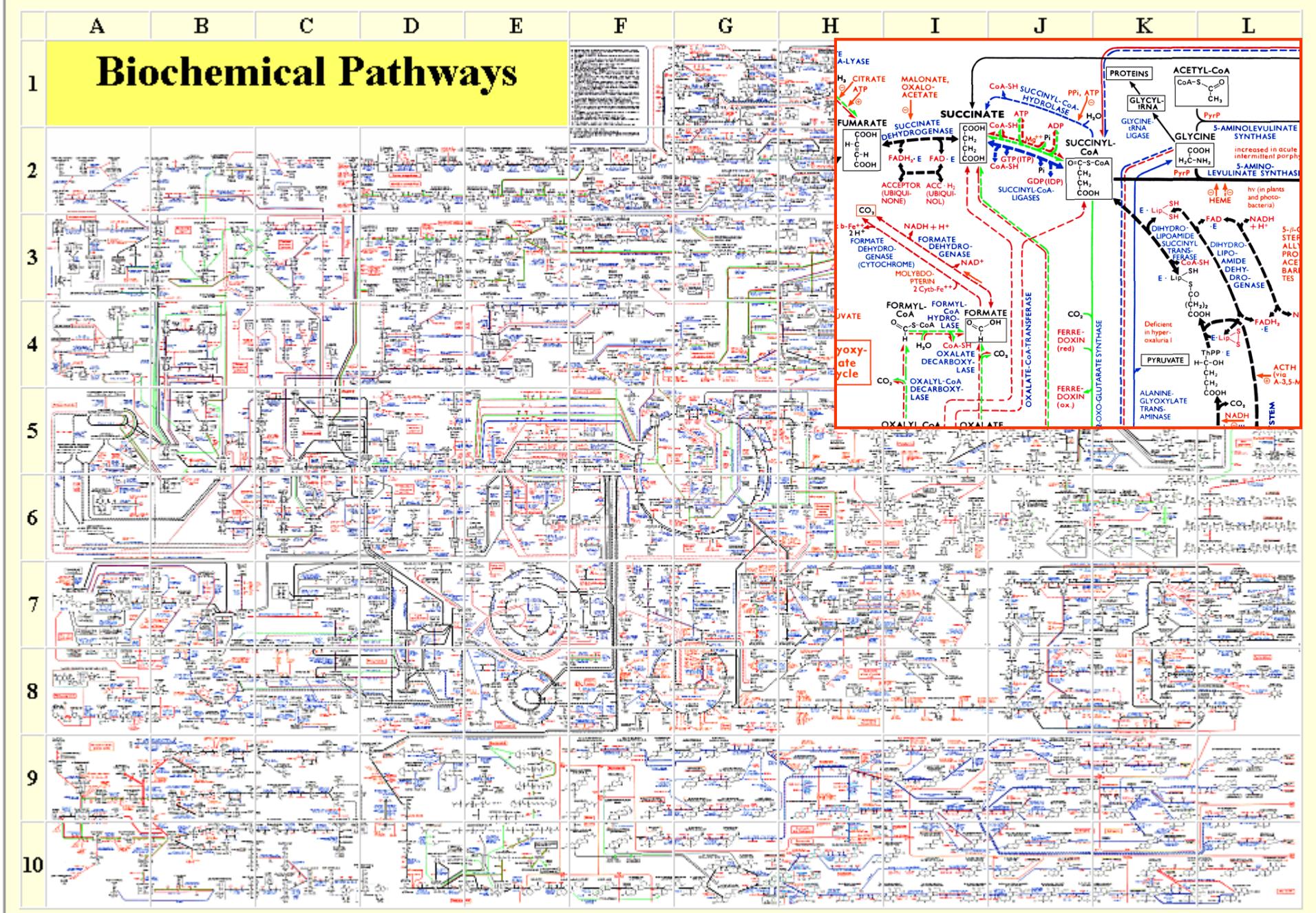


An induced subgraph of the collaboration graph (with Erdős number at most 2).

Made by Fan Chung Graham and Lincoln Lu in 2002.

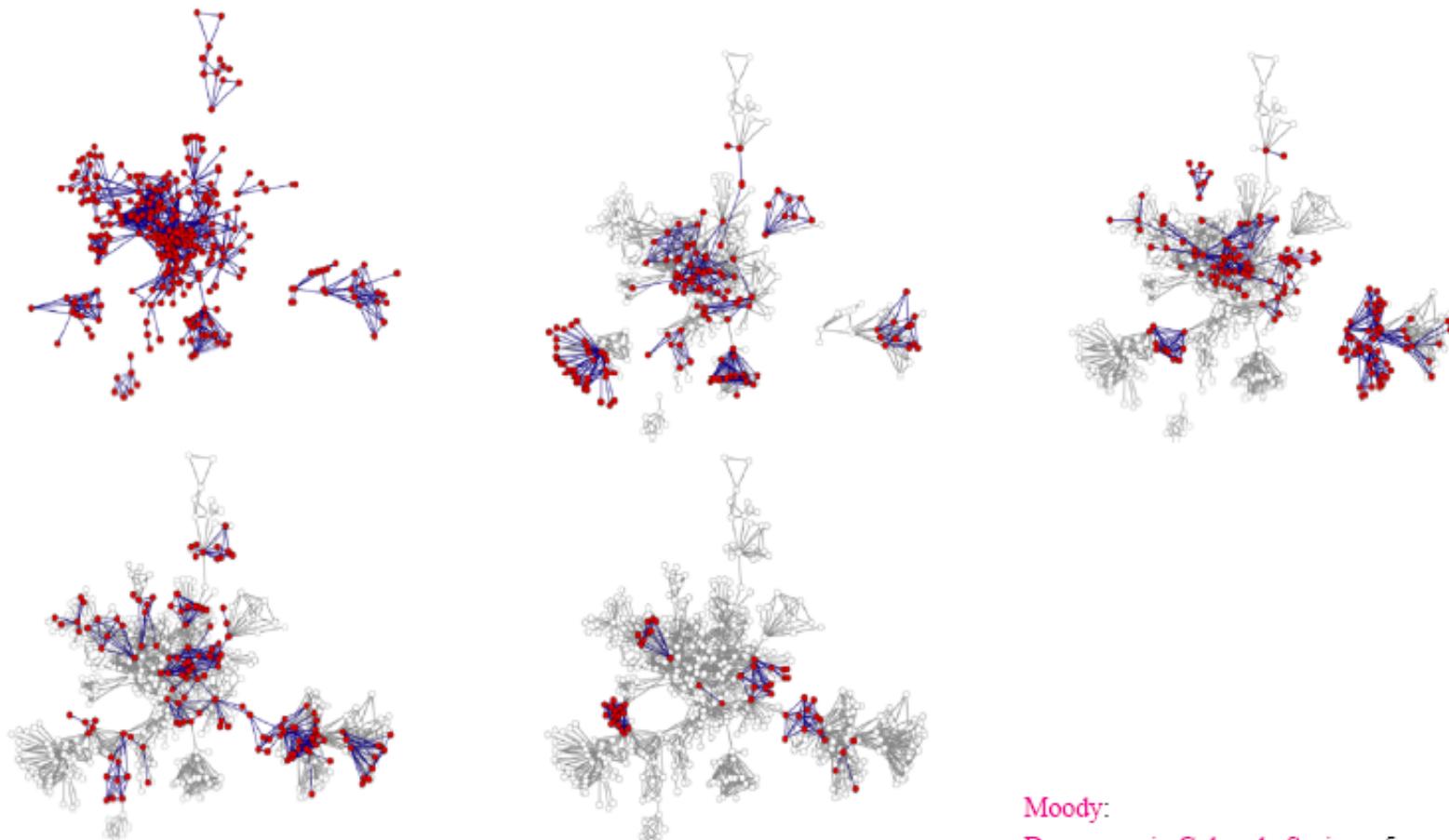
Erdos network: mathematician collaboration network (co-authorship network)

Biochemical Pathways



Temporal networks

In a *temporal network* the presence/activity of vertex/link can change through time. **Pajek** supports two types of descriptions of temporal networks based on *presence* and on *events*.



Moody:
Drug users in Colorado Springs, 5 years

Visualisation Challenge

1. Computational complexity

Efficiency

Runtime

We need more efficient algorithms

2. Visual complexity

Effectiveness

Readability

We need better untangling

- <http://www.visualcomplexity.com/vc/>

Main Approaches

Over the past 20 years, there have been several ideas put forward to solve the scale problem:

1. Cluster the data
2. Multi-level approach
3. Use 3 dimensions
4. Reduce Visual Complexity
5. Integration with Analysis
6. Integration with Interaction

1. Cluster the Data

Reduce size/dimension

- *Clustering*
- *Filtering*

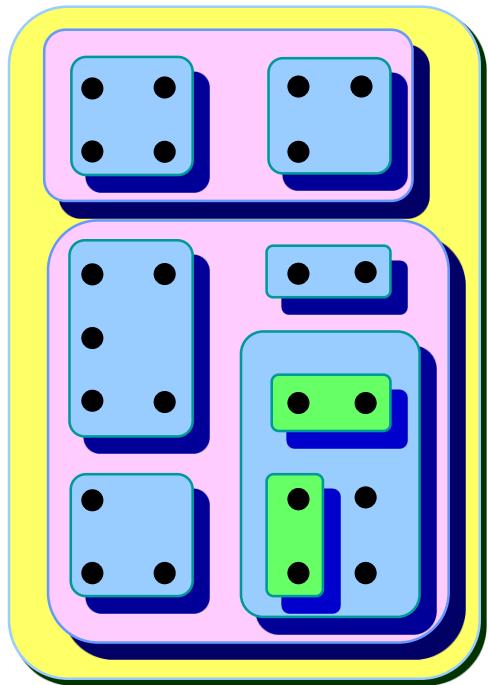
Clustering

There are many clustering algorithms

- ***Hierarchical clustering***
- K-means clustering
- Graph partitioning: Balanced min-cut
- Spectral clustering
- Community detection
- Infomap
- Label propagation
- Louvain

Metrics:

- Modularity
- Coverage
- ...



Clustered Graph Model

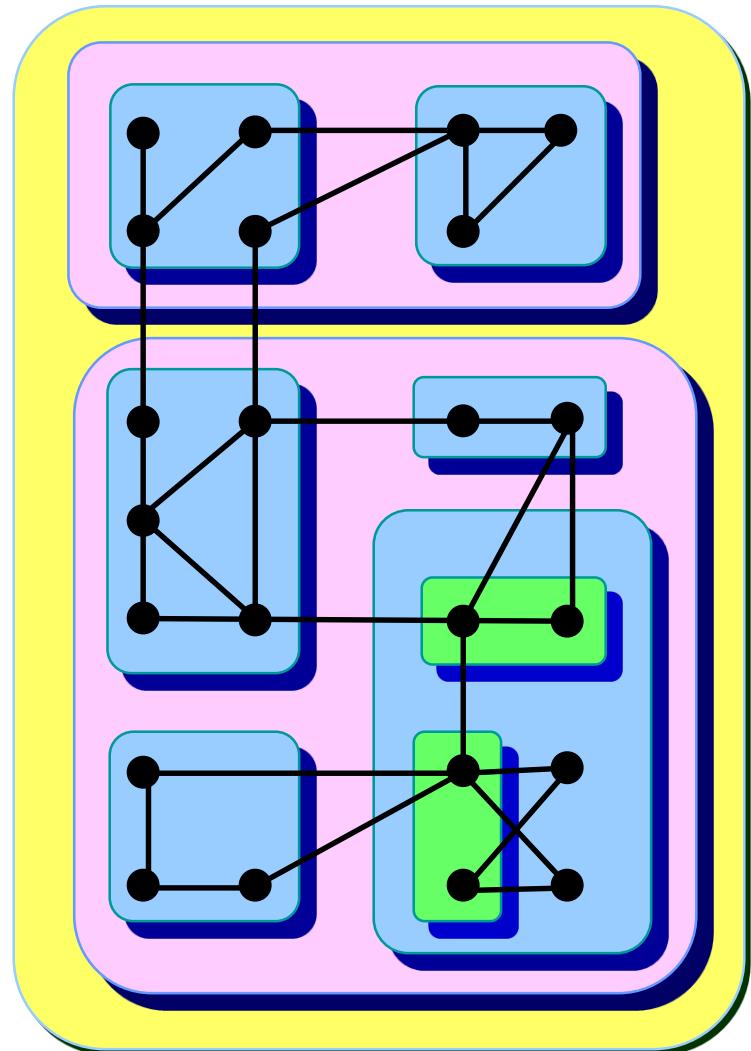
[Feng,Cohen,Eades 95]

A clustered graph $\mathbf{C}=(G,T)$ consists of

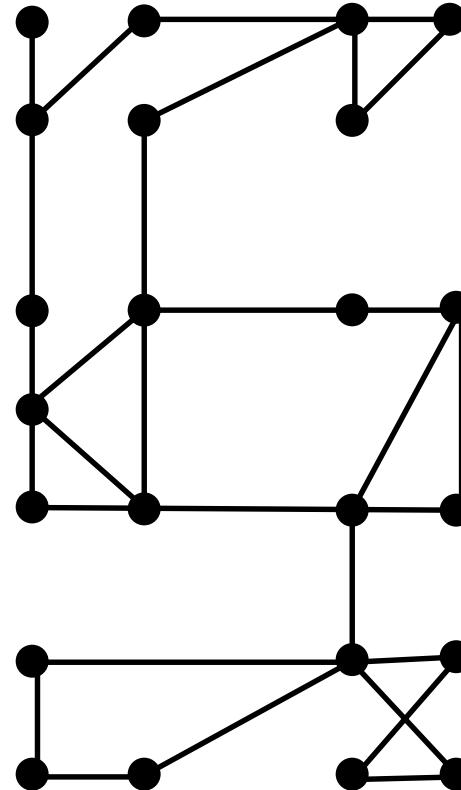
- a graph G , and
- a tree T

such that the leaves of the tree T are the vertices of G .

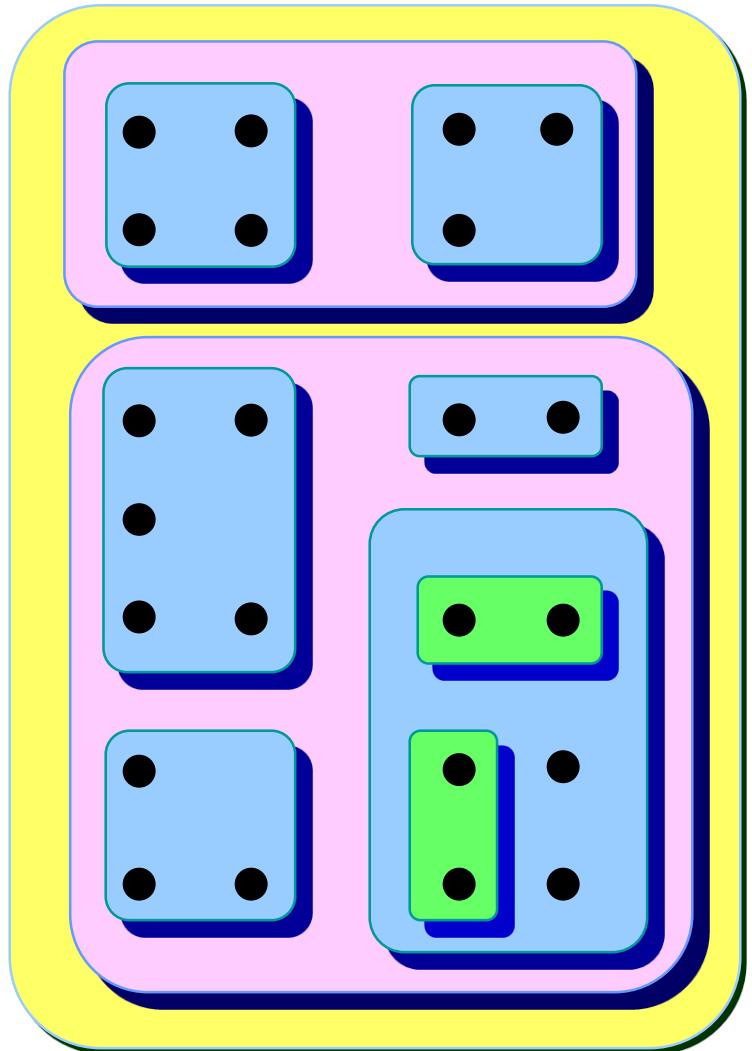
The tree T defines a clustering of the vertices of G .



graph

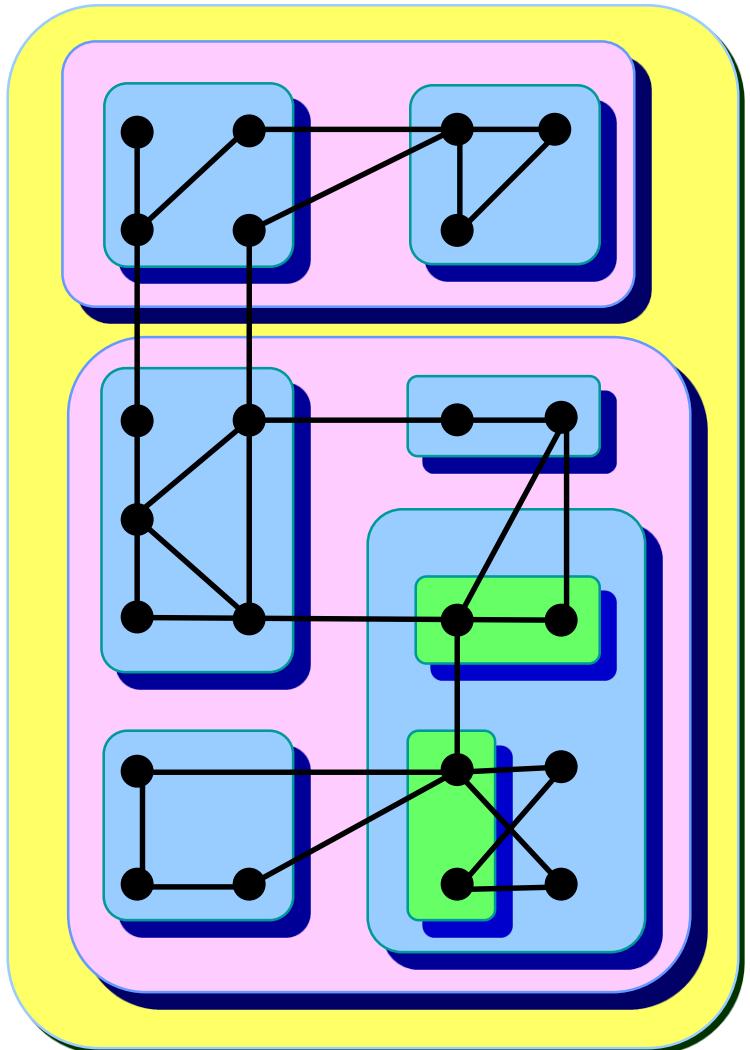


tree



graph
+
tree
=

clustered graph



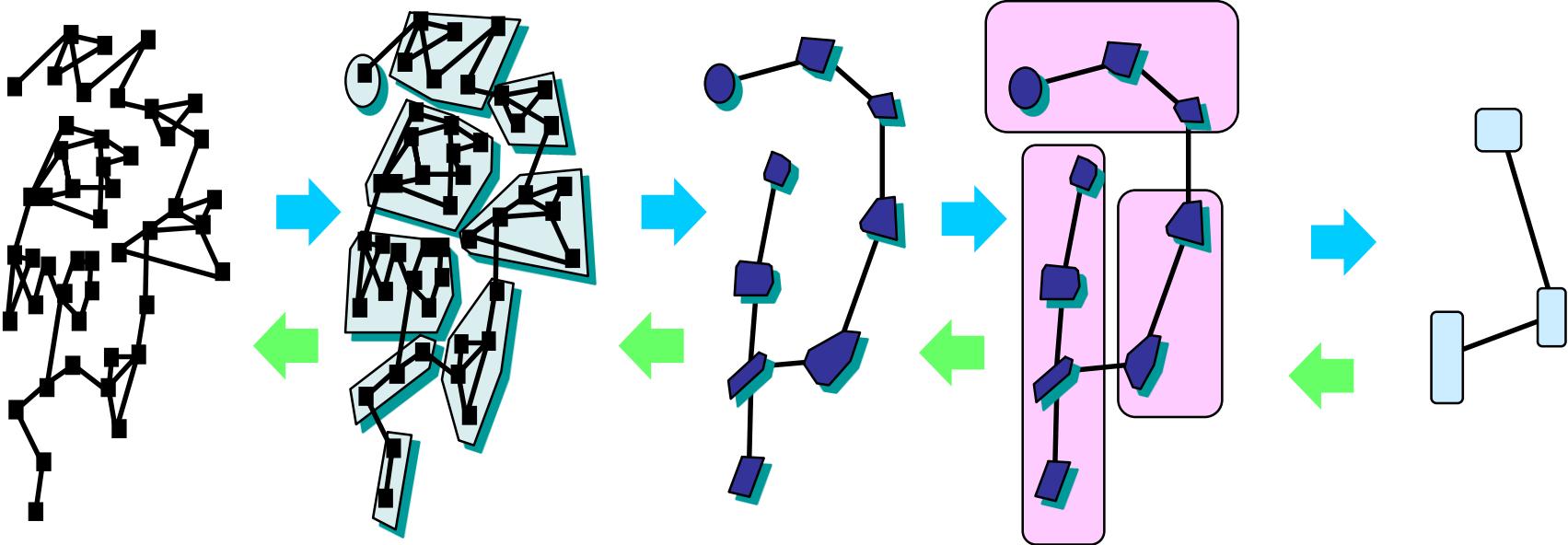
Clustered Graph Visualisation

Criteria:

- Avoid overlap between clusters
- Minimise edge crossing
- Minimise edge-cluster crossings

Simple solution:

- Weighted Spring algorithm with cluster constraints



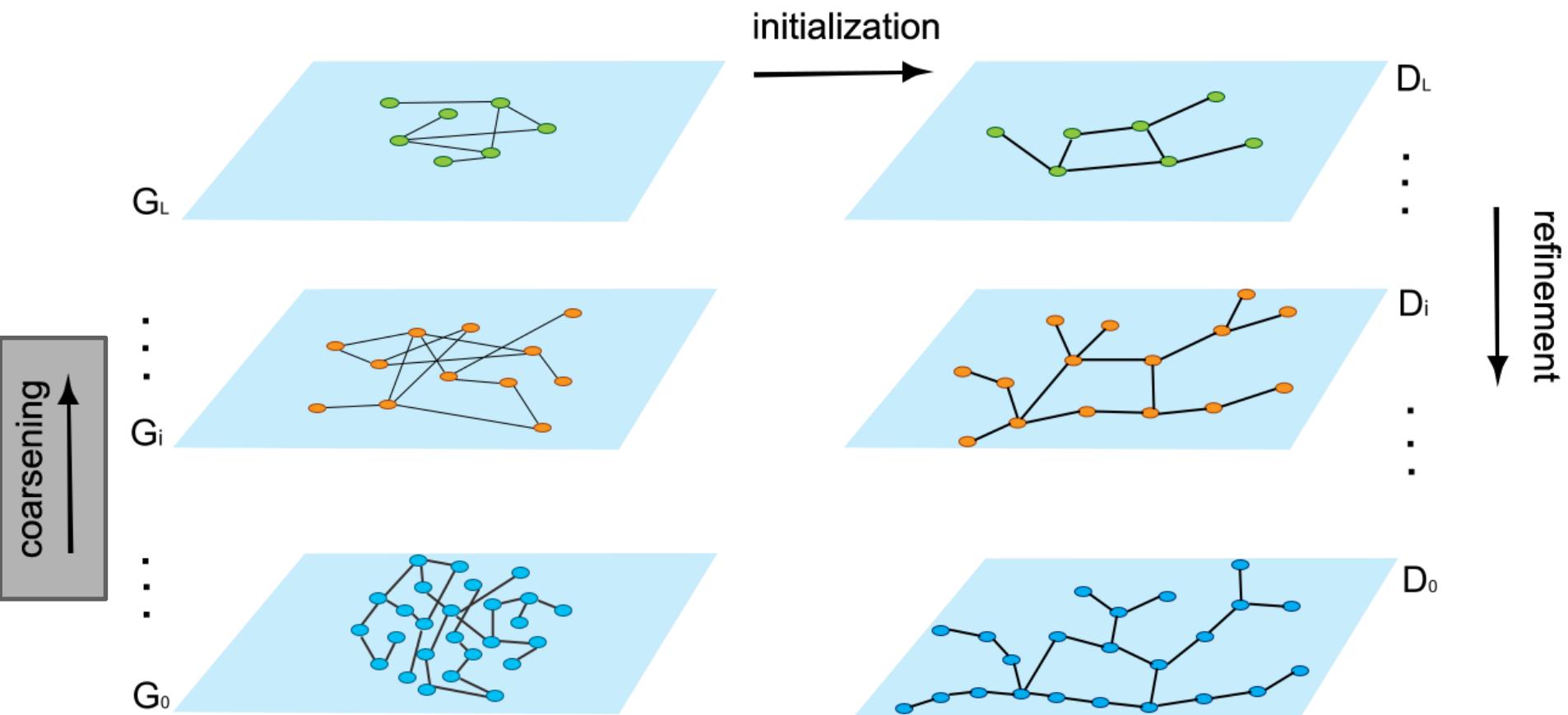
2. Multi-level Approach

Step 1: coarsening

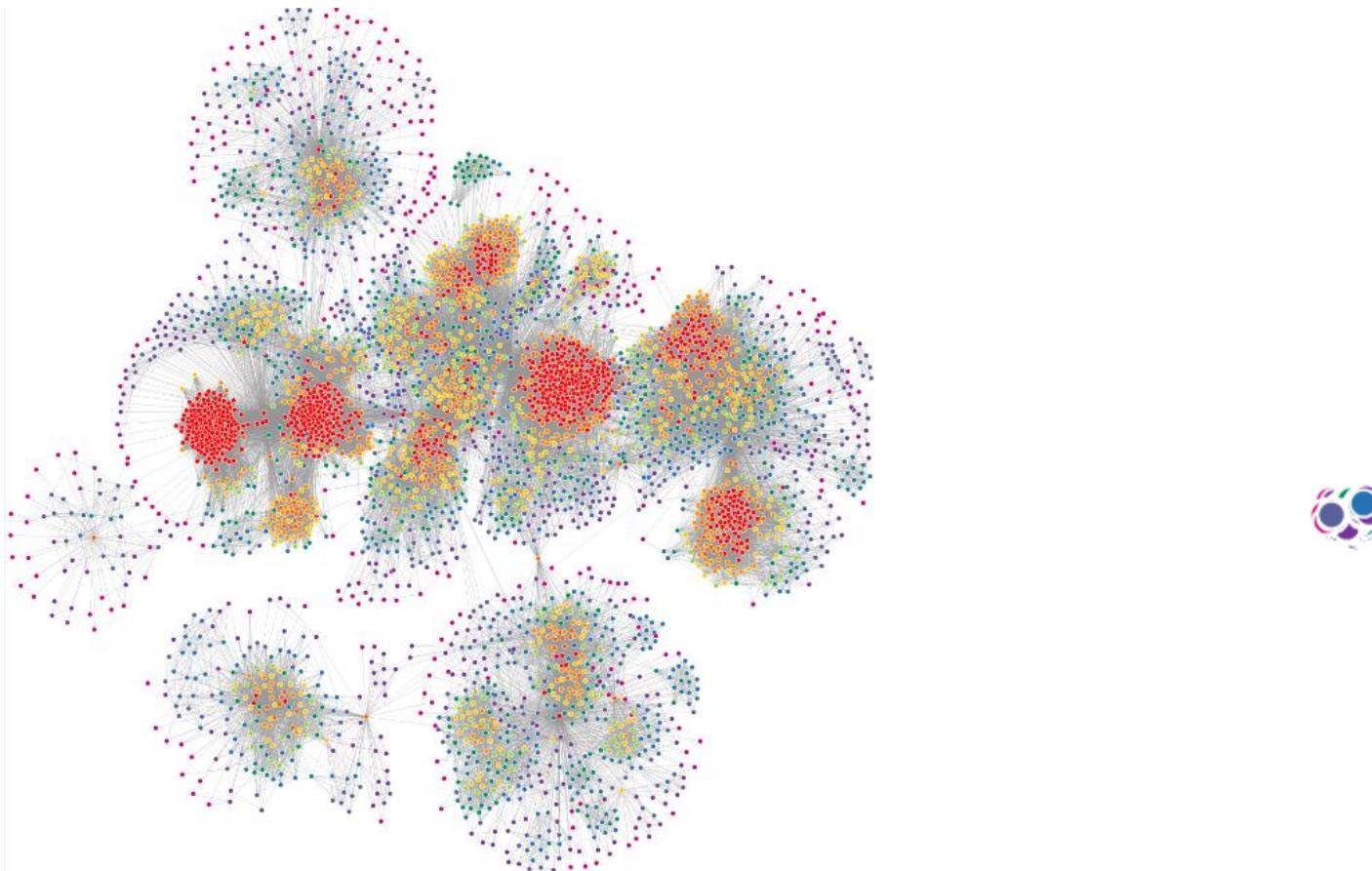
Step 2: refinement

Step 1: Coarsening Step

Step 2: Refinement Step



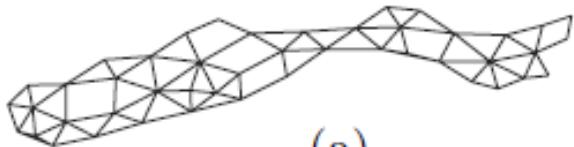
Quality (Untangling)



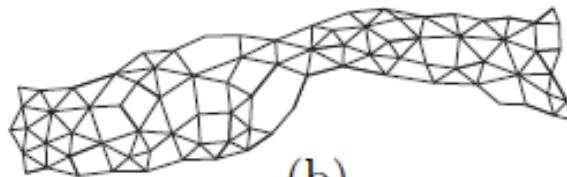
small Facebook with 4,039 nodes

[Chris Walshaw 00]

combinatorial clustering: graph partitioning



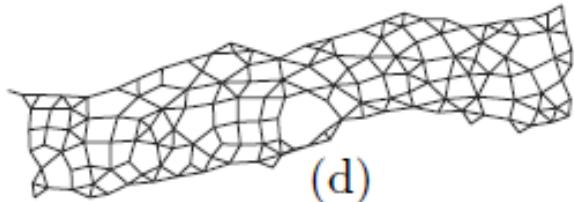
(a)



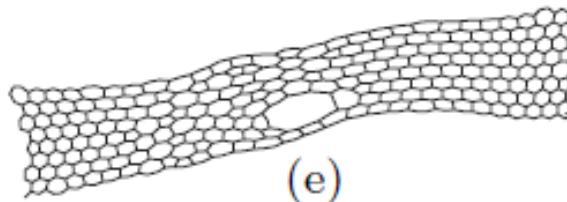
(b)



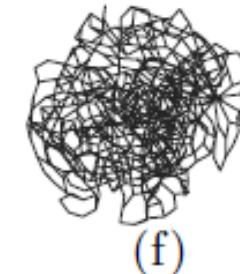
(c)



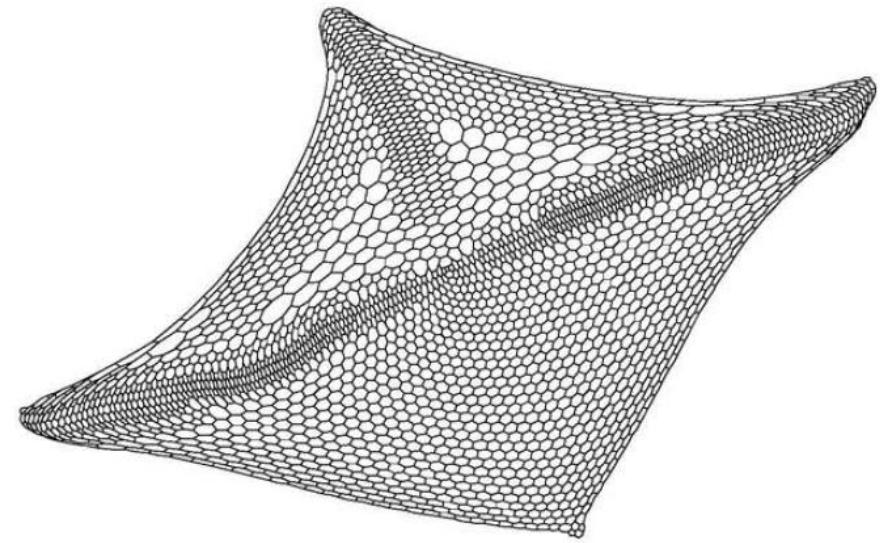
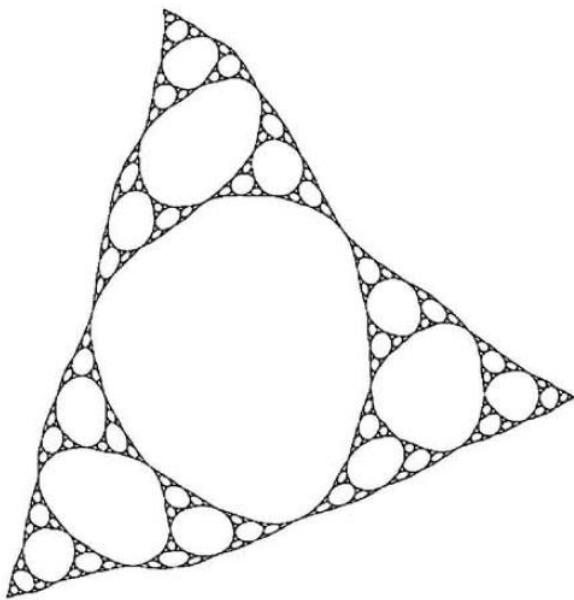
(d)



(e)

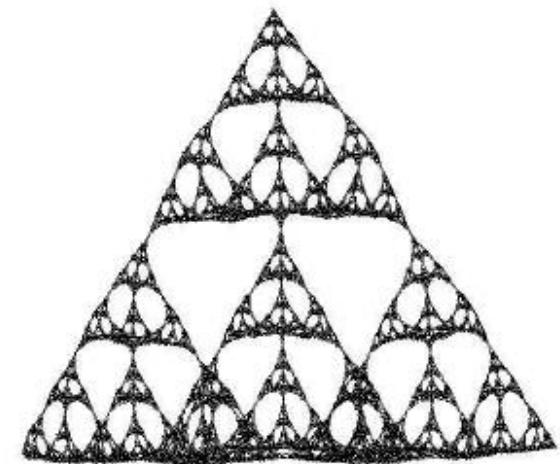
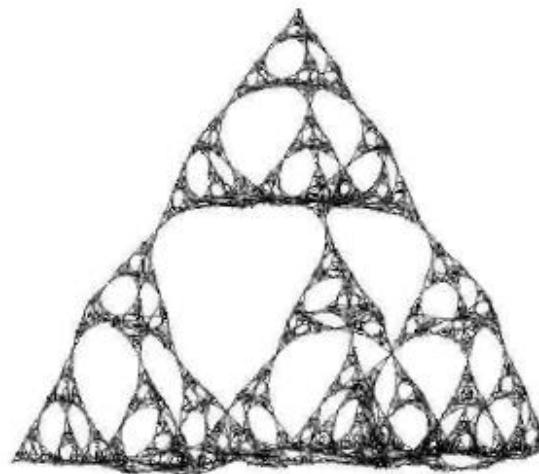
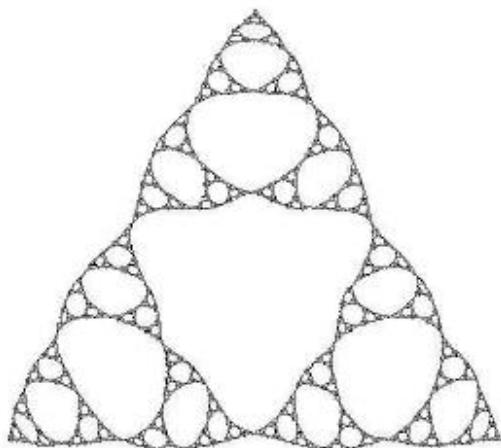
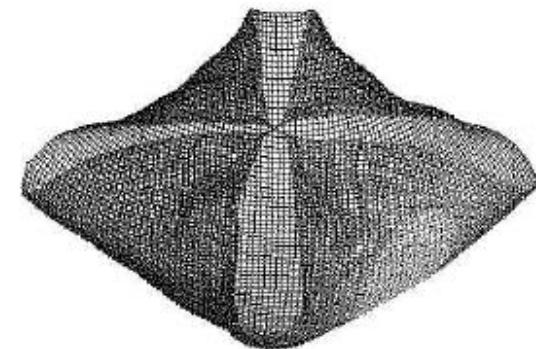
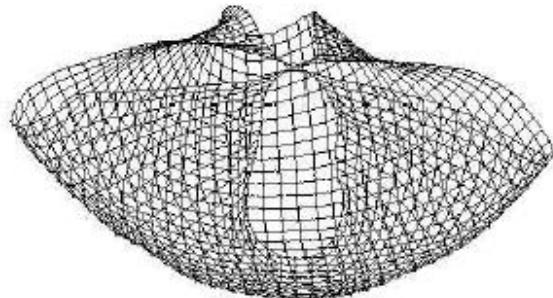


(f)



GRIP [Kobourov et al. 00]

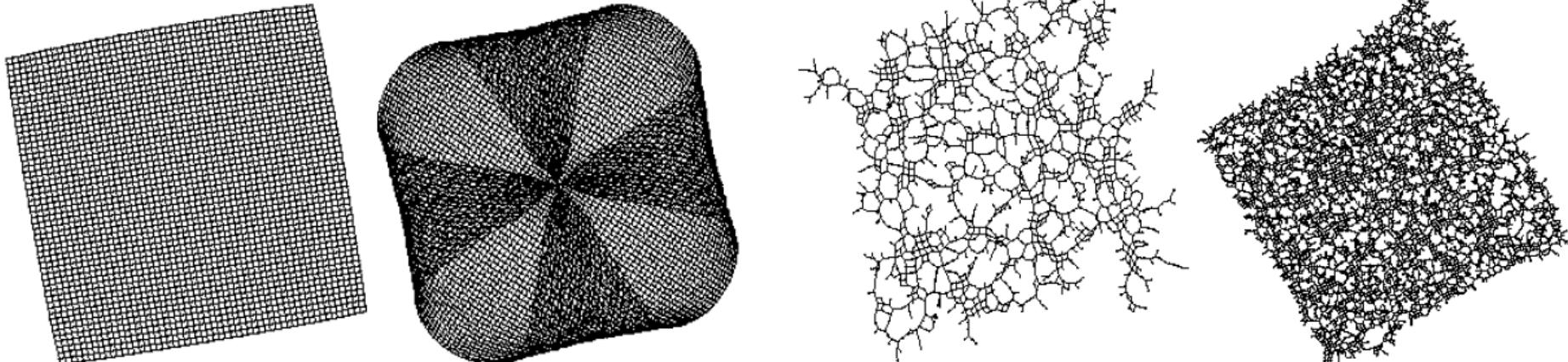
<https://www2.cs.arizona.edu/~kobourov/GRIP/>



Multi-scale Algorithm [Yehuda Koren 00]

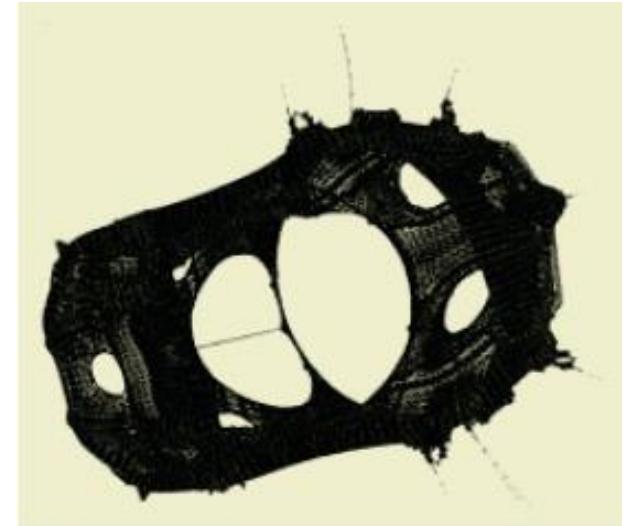
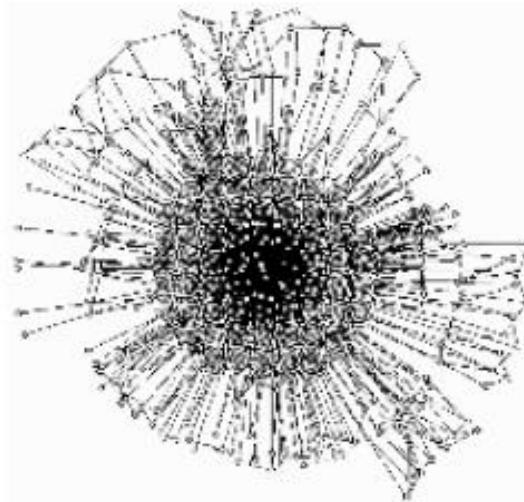
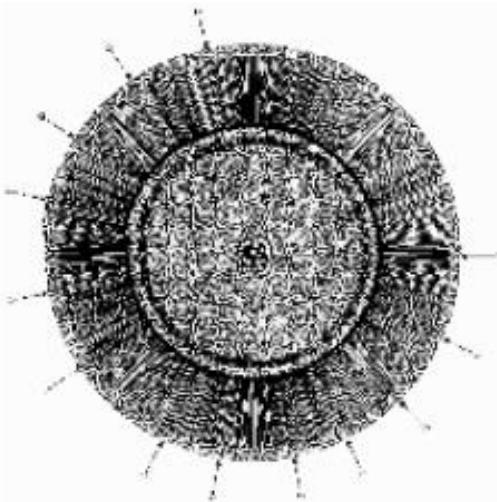
Other Methods by Koren:

- **HDE** [2002]: high dimensional embedding
- **ACE** [2002]: directed graph
- Spectral GD [2003]: spectral graph theory
- **DIG-COLA** [2005]: directed graph



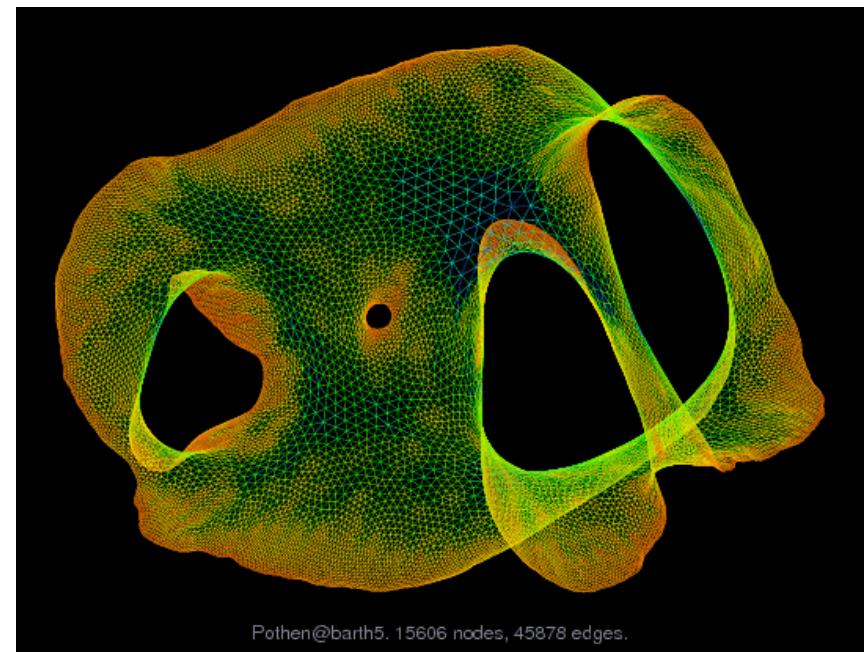
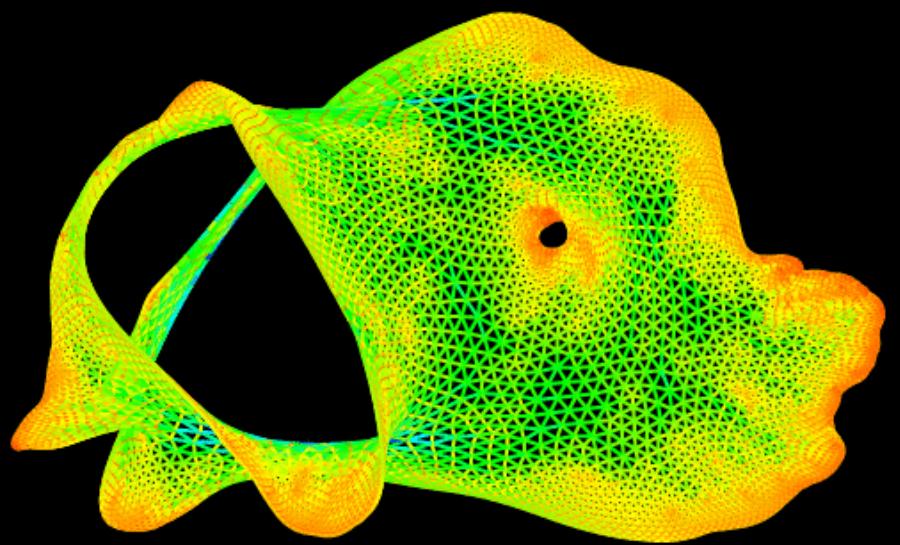
FM3 [Hachul & Junger 05]

- Similar to **FADE** algorithm: Solar merger using **Barnes-Hut** method (*approximate n-body simulation problem in $O(n\log n)$ time using the Quad/Oct tree*)
- Implementation: OGDF
- <http://amber-v7.cs.tu-dortmund.de/doku.php/start>



sfdp [Yifan Hu]

<http://yifanhu.net/SOFTWARE/SFDP/index.html>
Graphviz



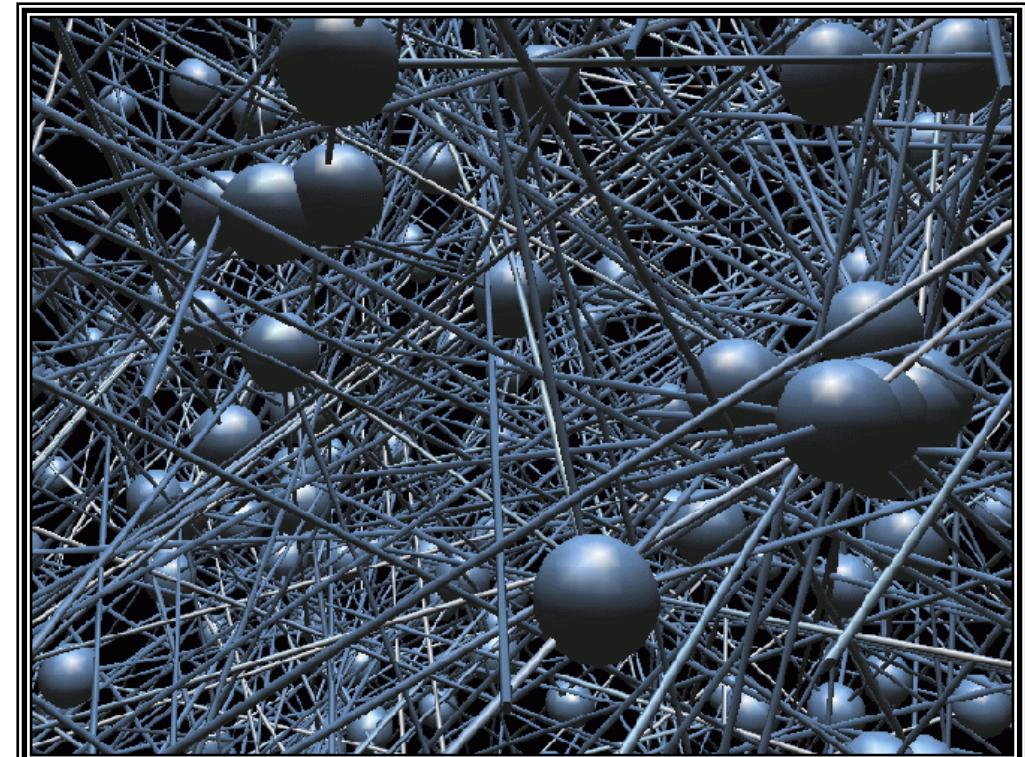
3. Use the 3rd Dimension

2.5D Network Visualisation Framework

- *Trees*
- *Clustered graphs*
- *Directed/Hierarchical graphs*
- *Dynamic/Temporal networks*
- *Multi-relational networks*

3D Graph Visualisation

- 1990 – 2005:
 - Many theorems on 3D
 - Many metaphors
 - Many research grants
 - Many experiments with 3D network visualization
 - A start-up company

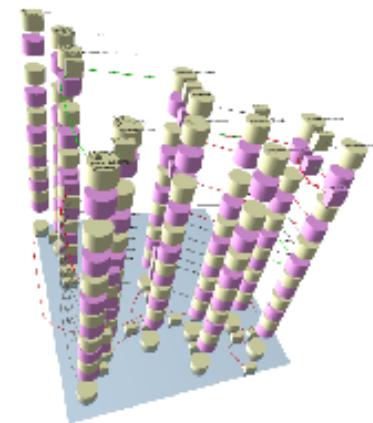
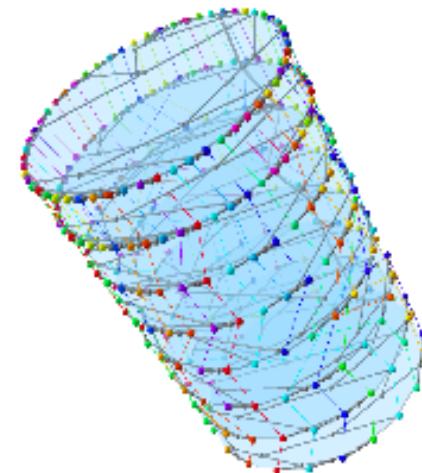
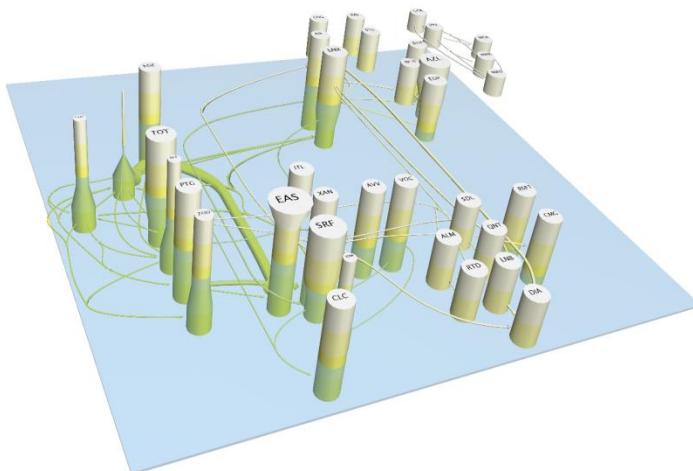


Mostly, 3D failed.

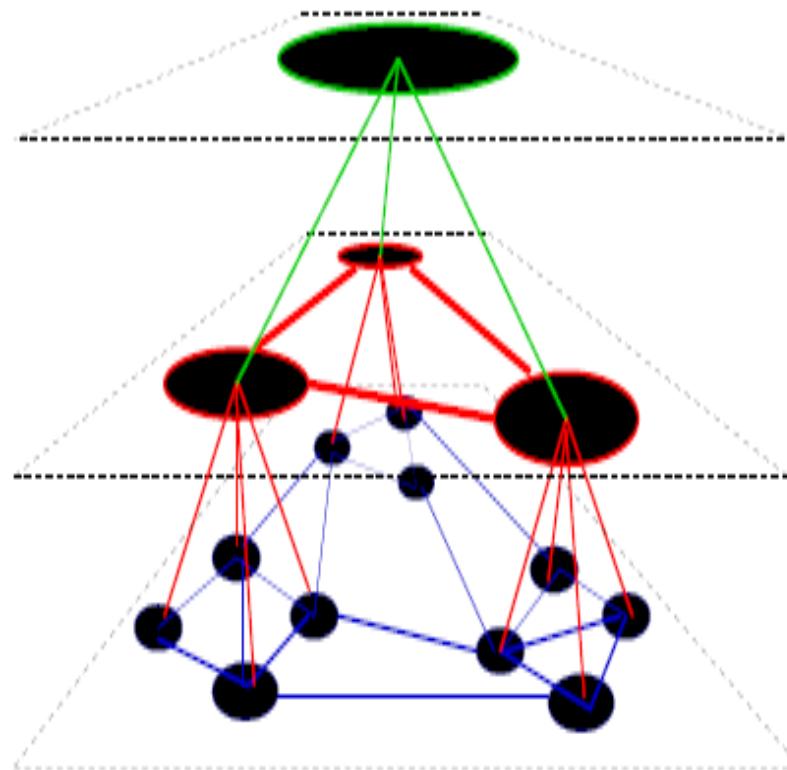
2.5D Graph Visualisation

2003+: success with 2.5D

- Colin Ware: “use 3D with a 2D attitude”
- Brandes/Dwyer 2003: 2.5D for temporal network
- Hong et al. 2003-6: *Multiplane framework*

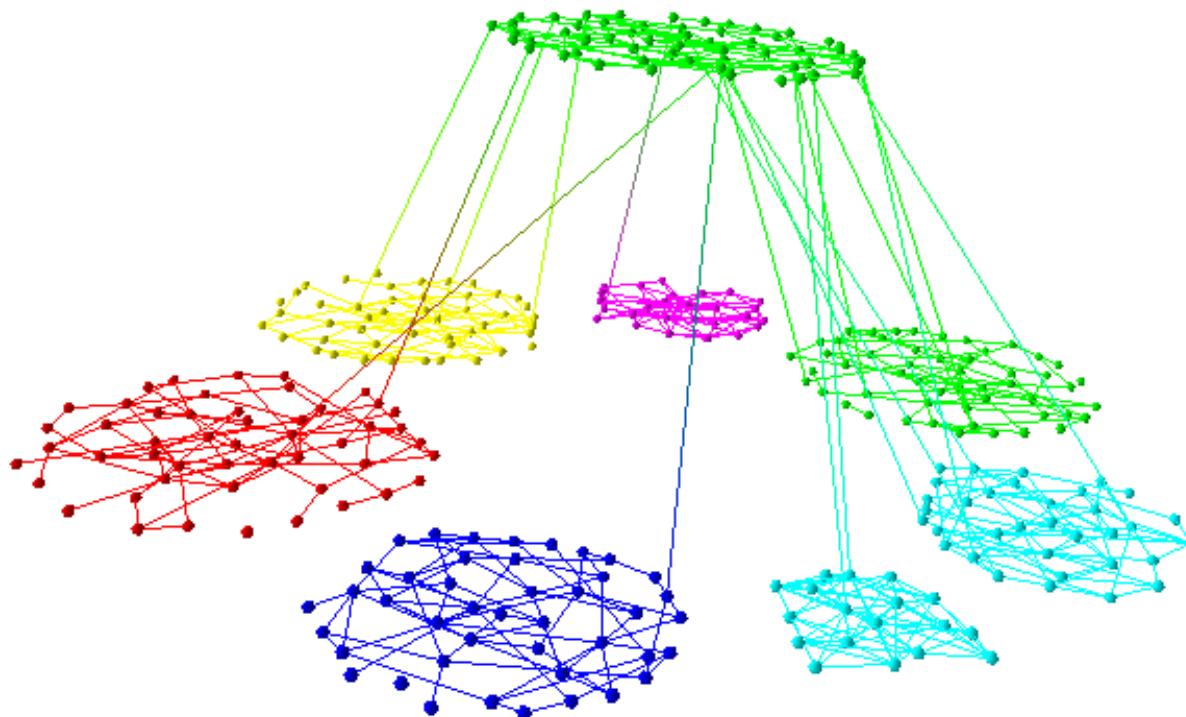


Multi-level Representation of Clustered Graph [Feng, Eades 96]

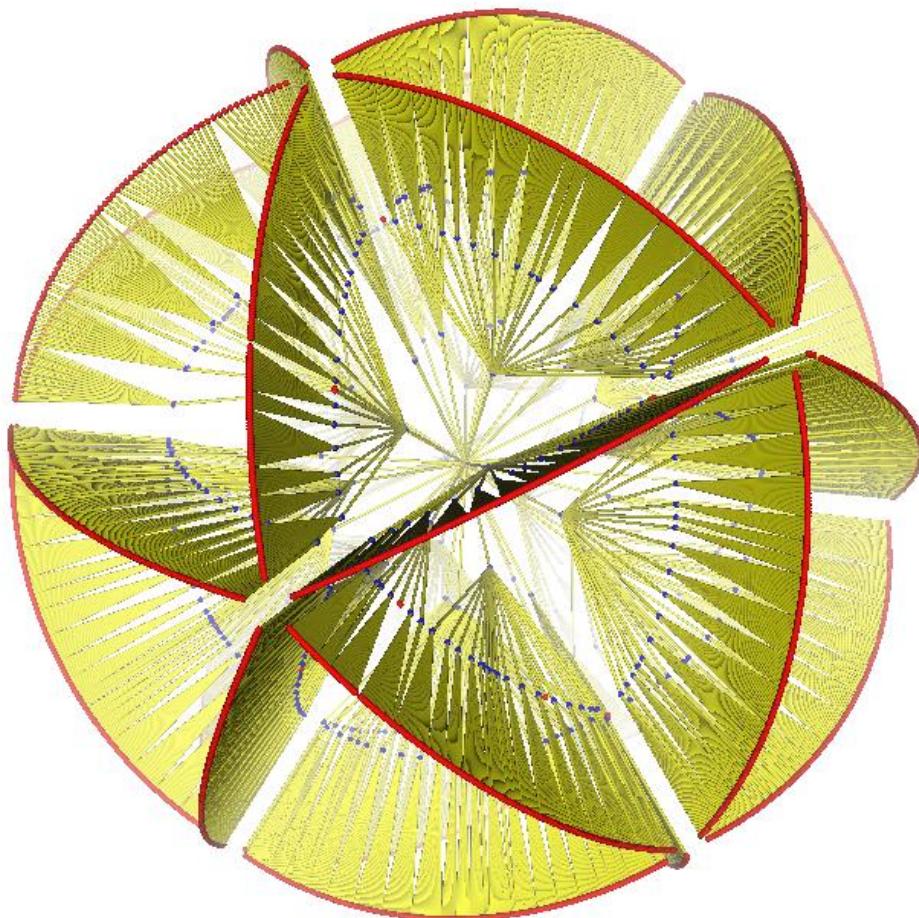


MultiPlane Algorithm [Hong 05]

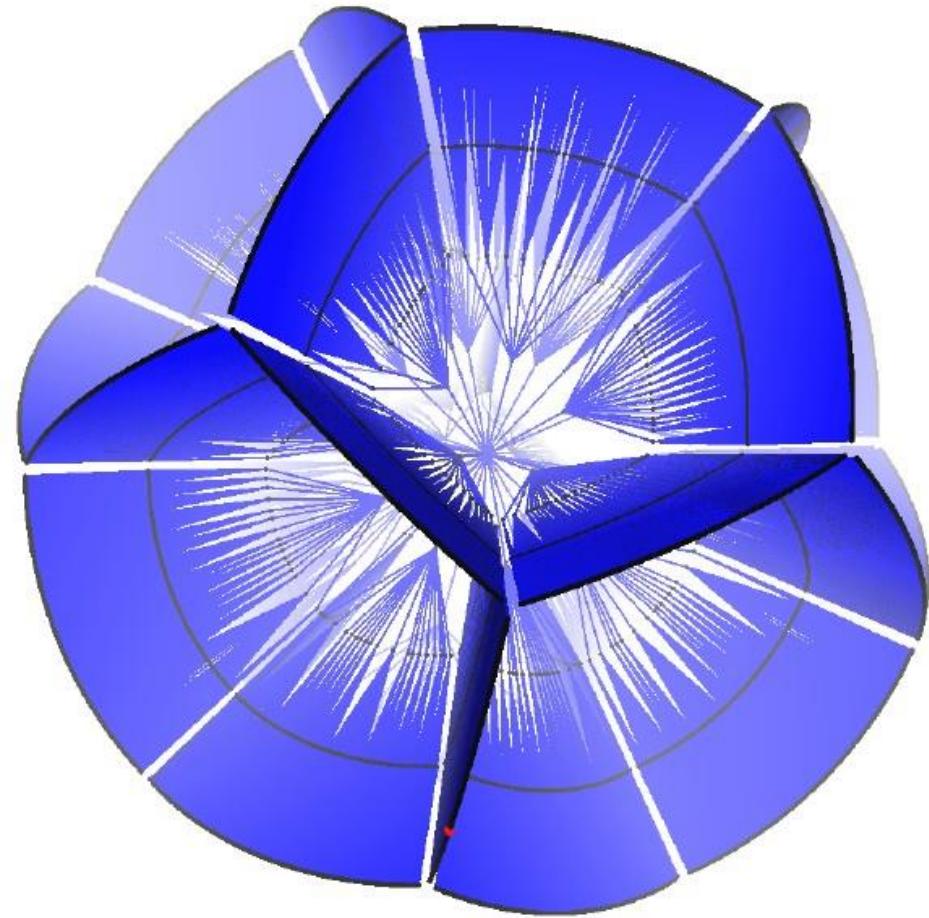
1. Divide the large/complex network into clusters
2. Draw each cluster on a 2D manifold in 3D
3. Connect the planes with inter-plane edges



2.5D Tree Layout: PolyPlane

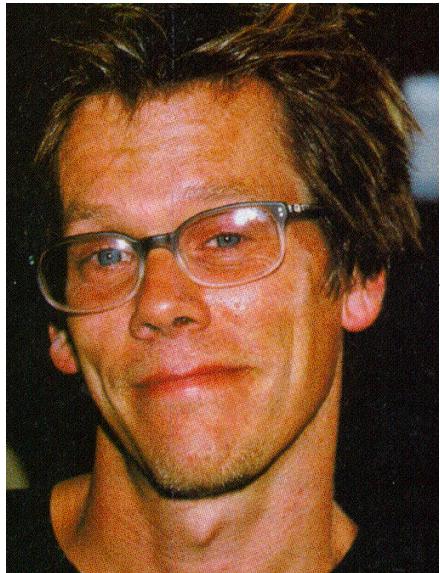


[Hong, Murtagh 03]
 $|V|=6929$ (30 planes)

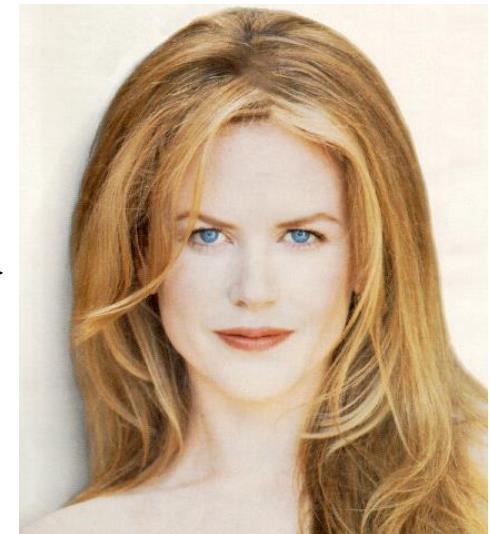
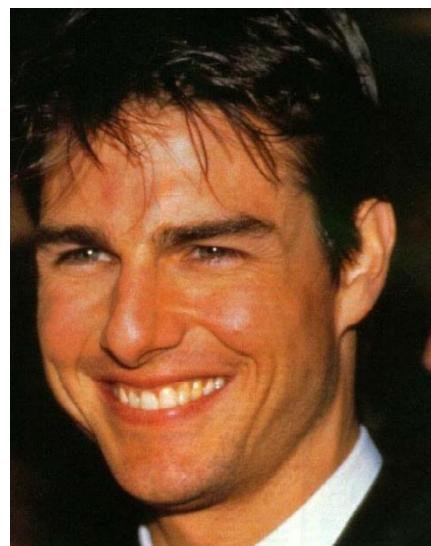


$|V|=146716$ (30 planes)

Hollywood Movie Actor Collaboration Network



A Few Good Men

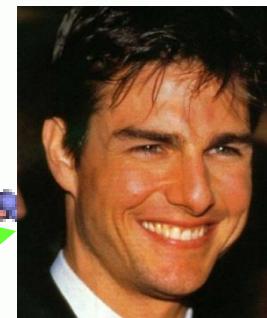
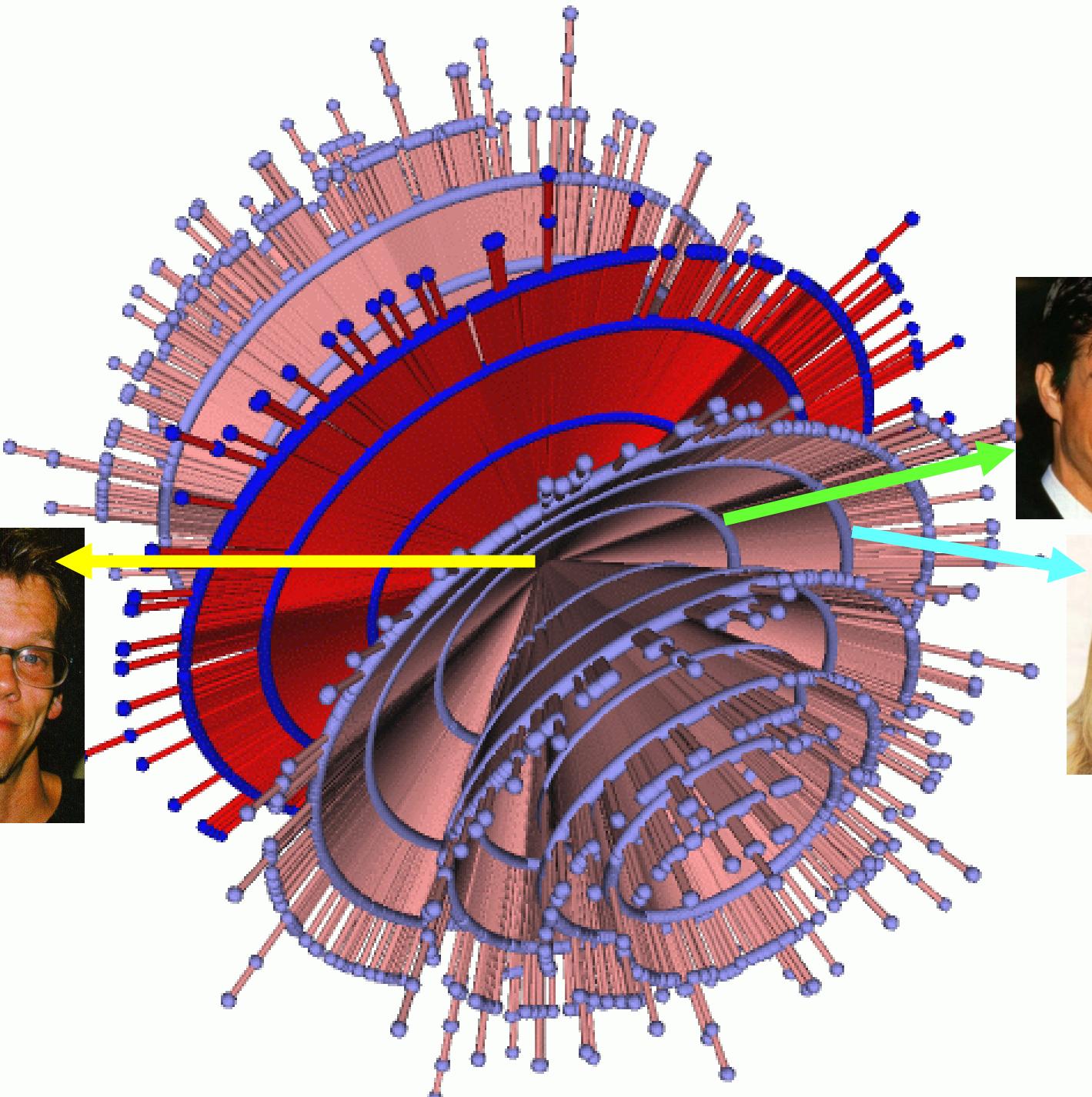
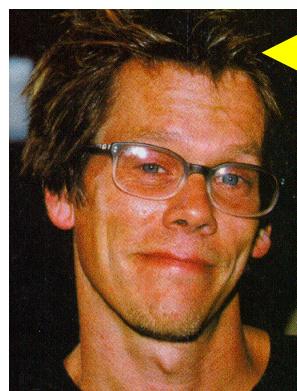


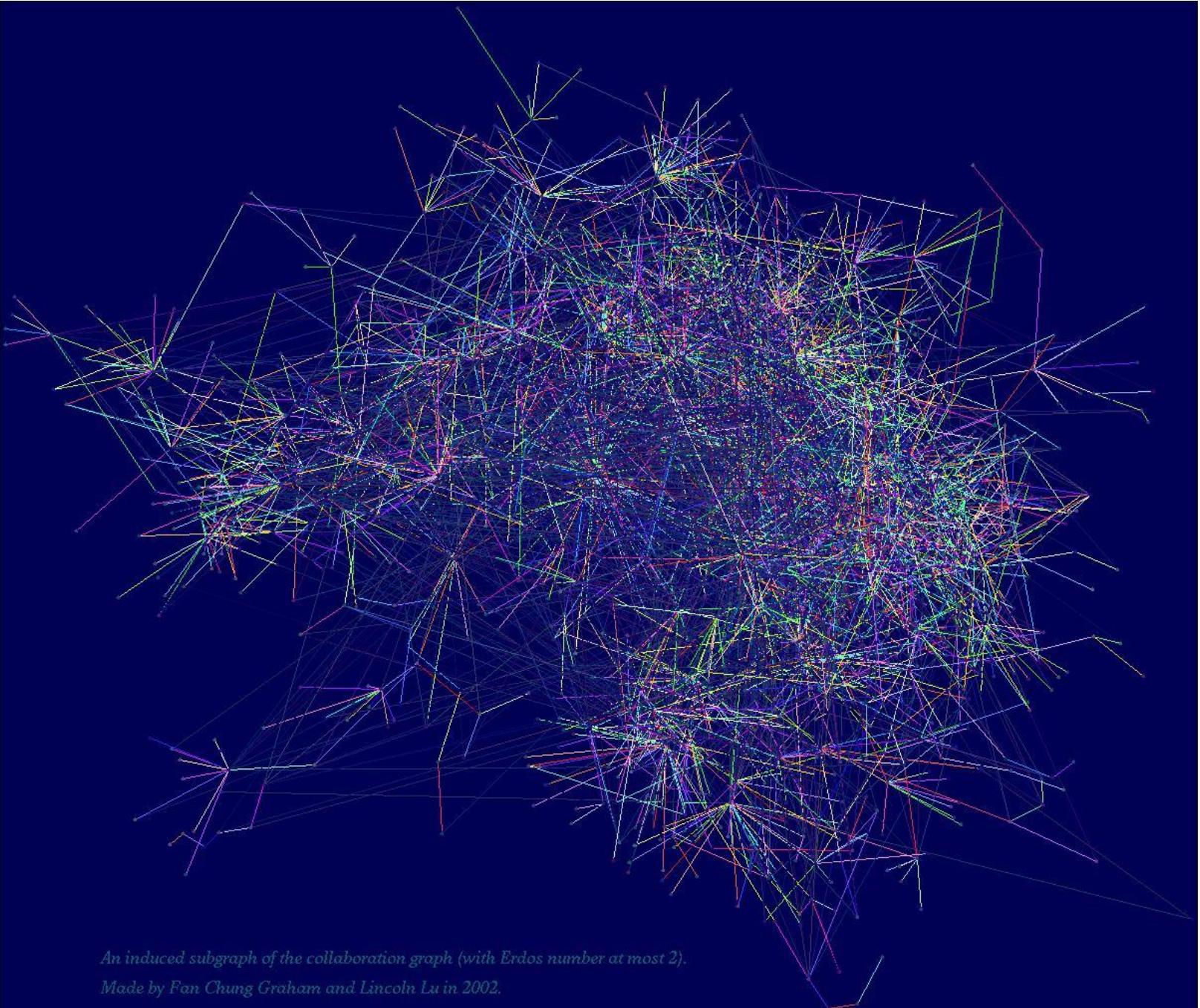
Days of Thunder (1990)

Far and Away (1992)

Actor Collaboration Network

Kevin Bacon Network

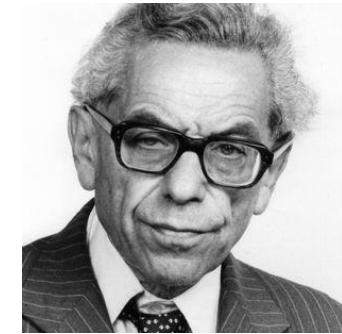
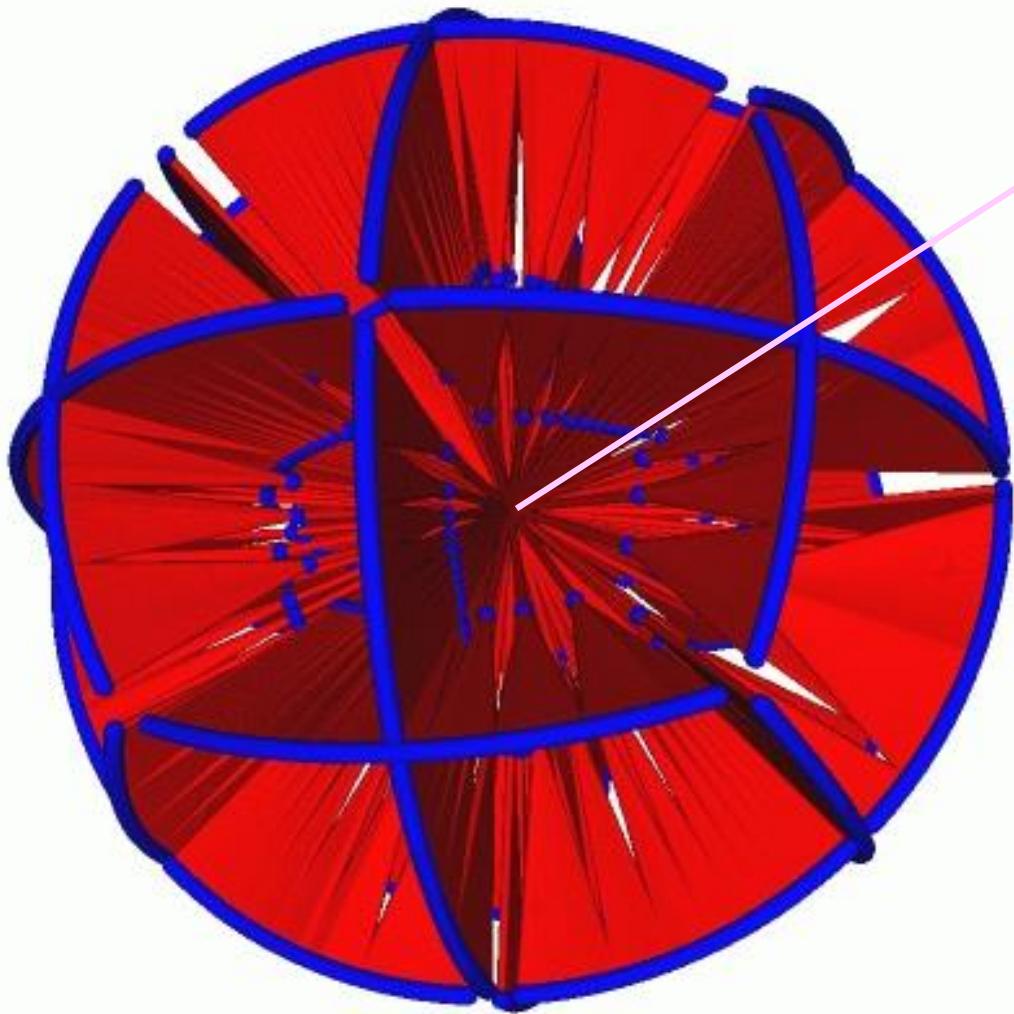




An induced subgraph of the collaboration graph (with Erdos number at most 2).

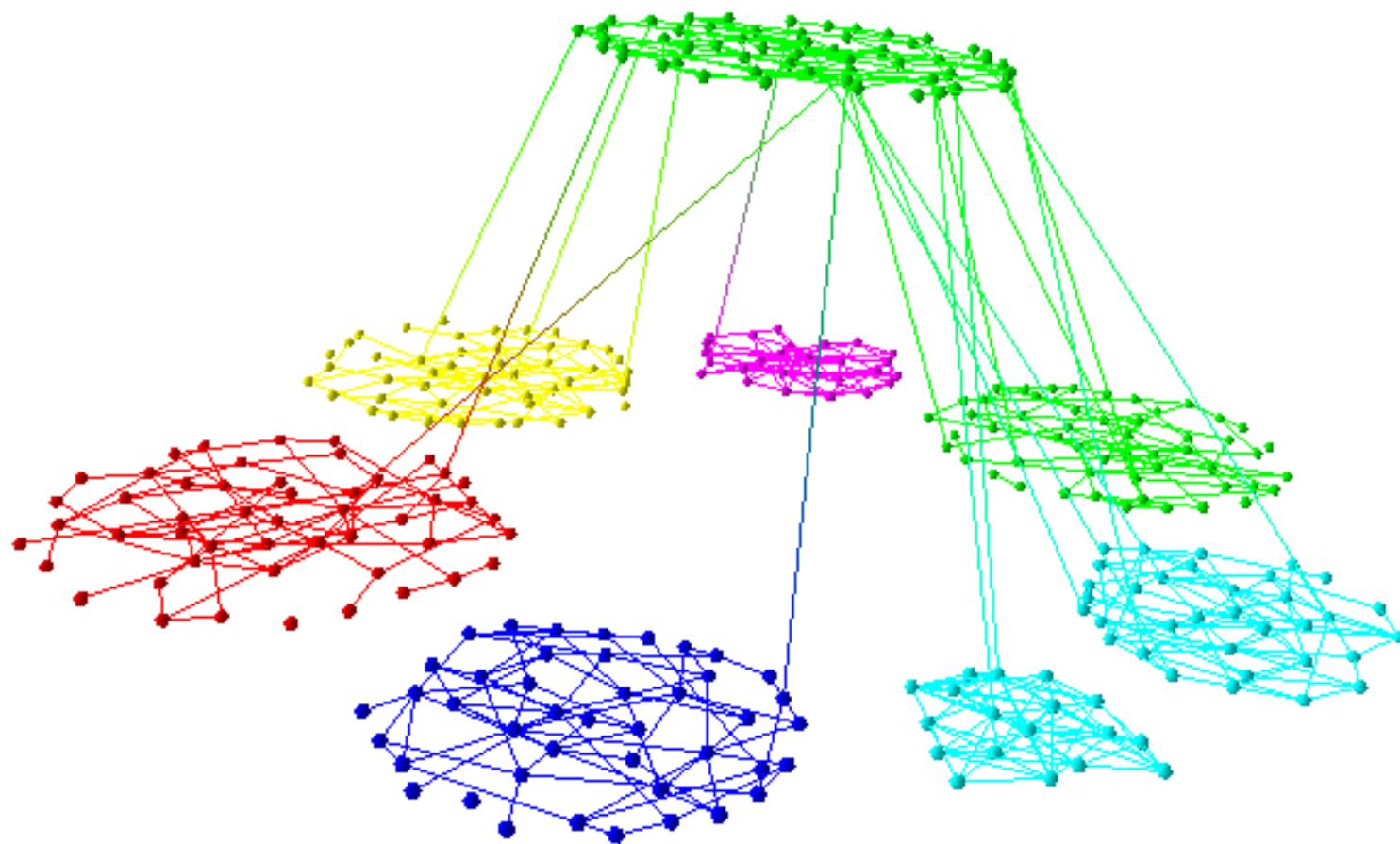
Made by Fan Chung Graham and Lincoln Lu in 2002.

Erdos network: mathematician collaboration network (co-authorship network)



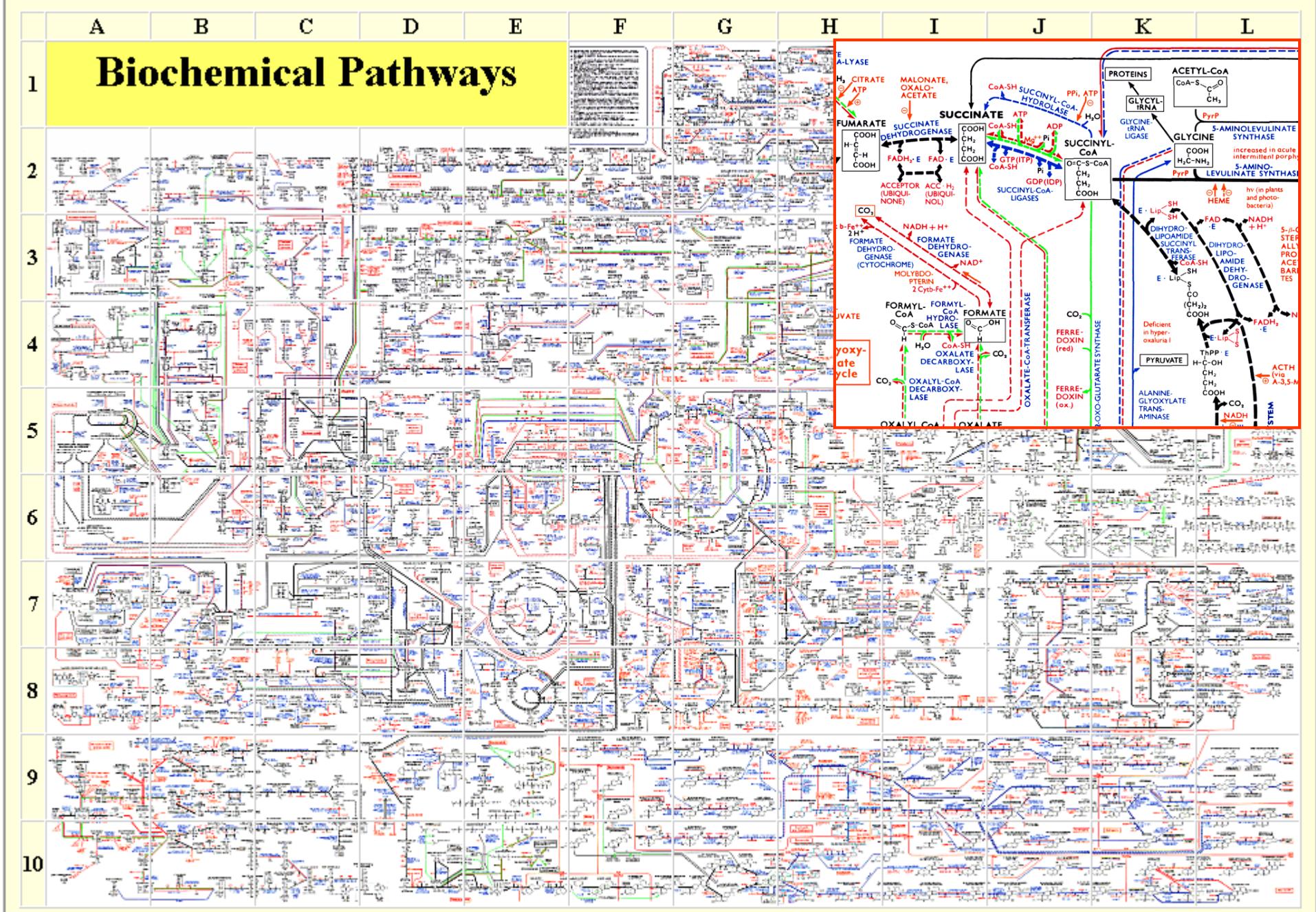
Erdős network: mathematician collaboration network (co-authorship network)

2.5D Clustered Graph Layout

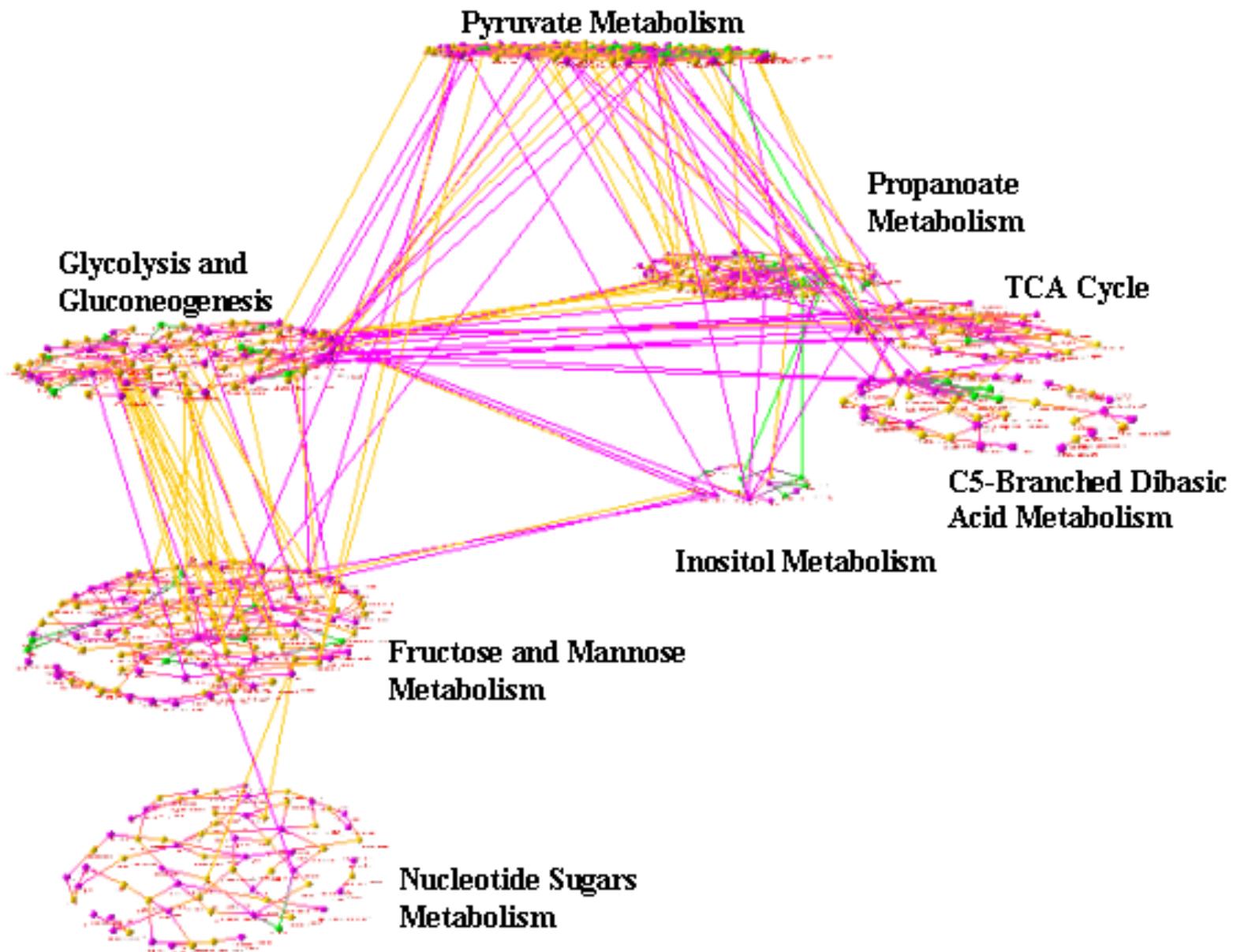


[Ho, Hong 05] Weighted 3D Tree Drawing Algorithm: $O(n)$ time
+ 2D Spring Algorithm: $O(n^2)$ time

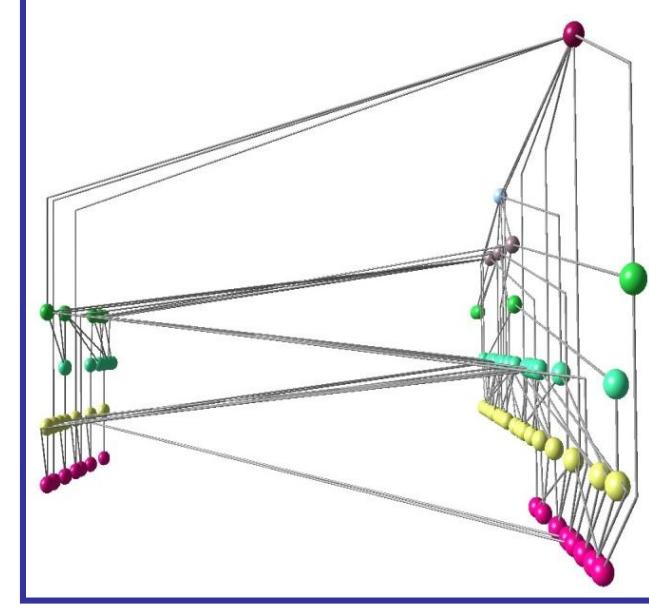
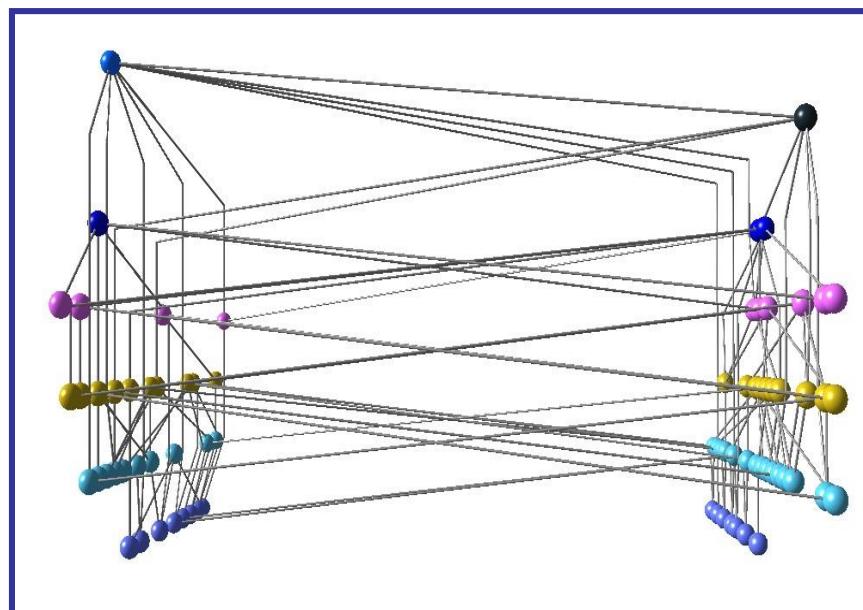
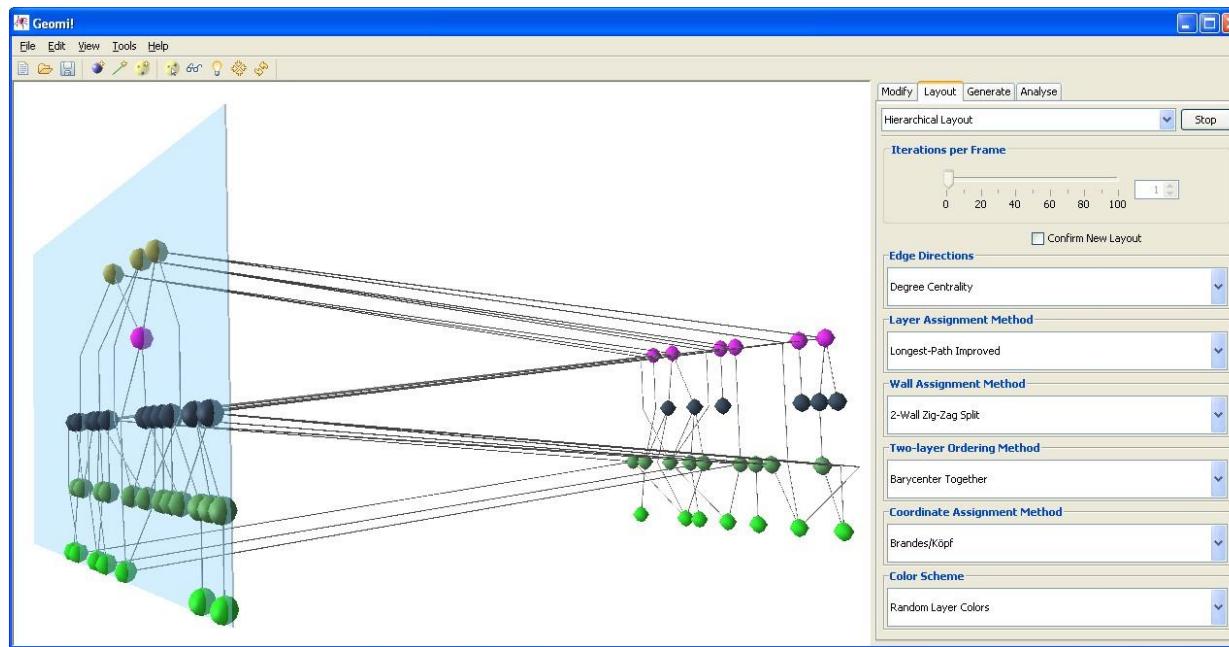
Biochemical Pathways



8 Related Metabolic Pathway Visualisation



2.5D Sugiyama Method [Nikolov, Hong 05/06]



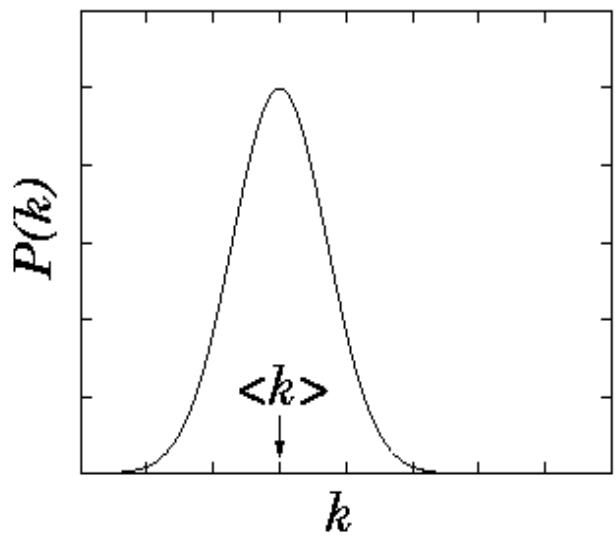
2.5D Scale-free Network Layout

[Ahmed et al. 05]

Scale-free Network

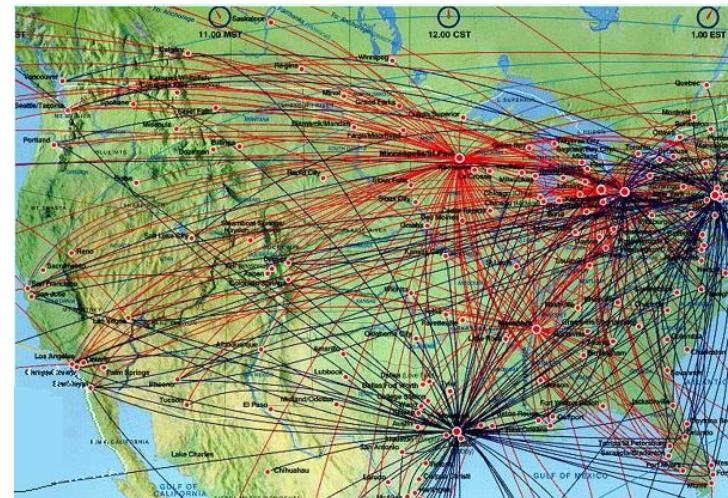
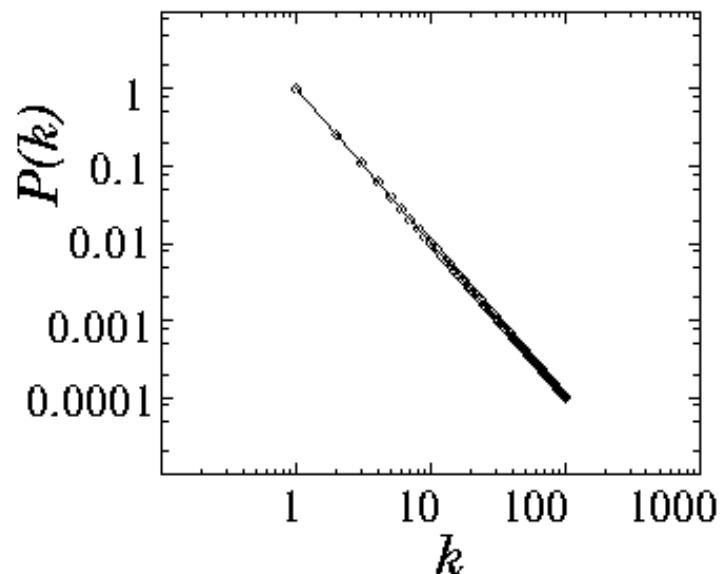
- [Barabasi and Albert 99]
 - Exponential Growth
 - Preferential attachment
 - Properties
 - Power-law degree distribution
 - Sparse, but locally dense
 - Small-world property: $O(\log \log n)$ average path length
 - High clustering coefficient
 - Examples:
 - Webgraph, Social networks, Biological networks
- [Barabasi和Albert 99]
- 指数增长
- 优惠附件
• 属性
- 幂律学位分布
- 稀疏, 但局部密集
- 小世界属性: $O(\log \log n)$ 平均路径长度 - 高聚类系数
• 例子:
- Webgraph, 社交网络, 生物网络

Poisson distribution

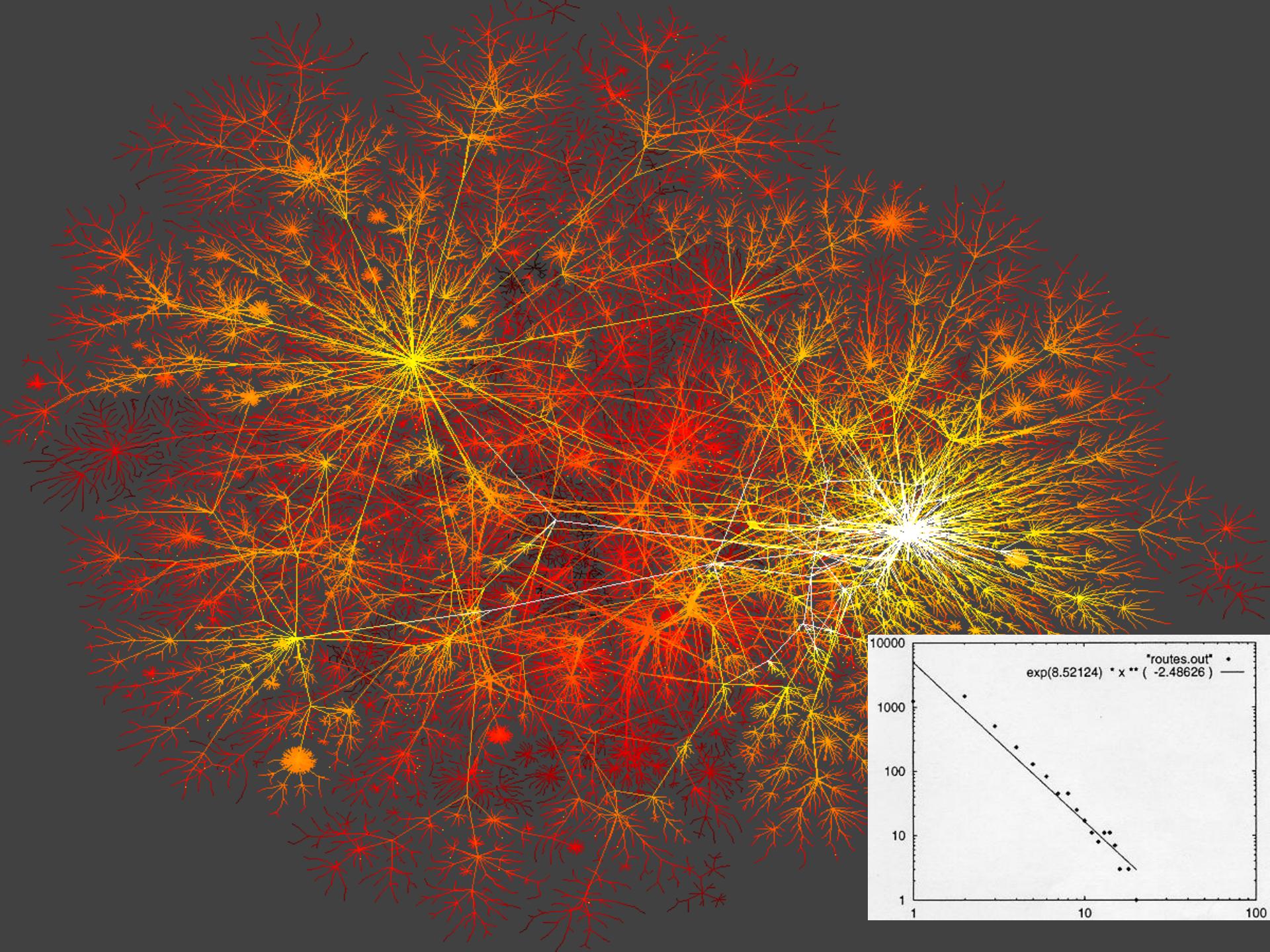


Exponential Network

Power-law distribution



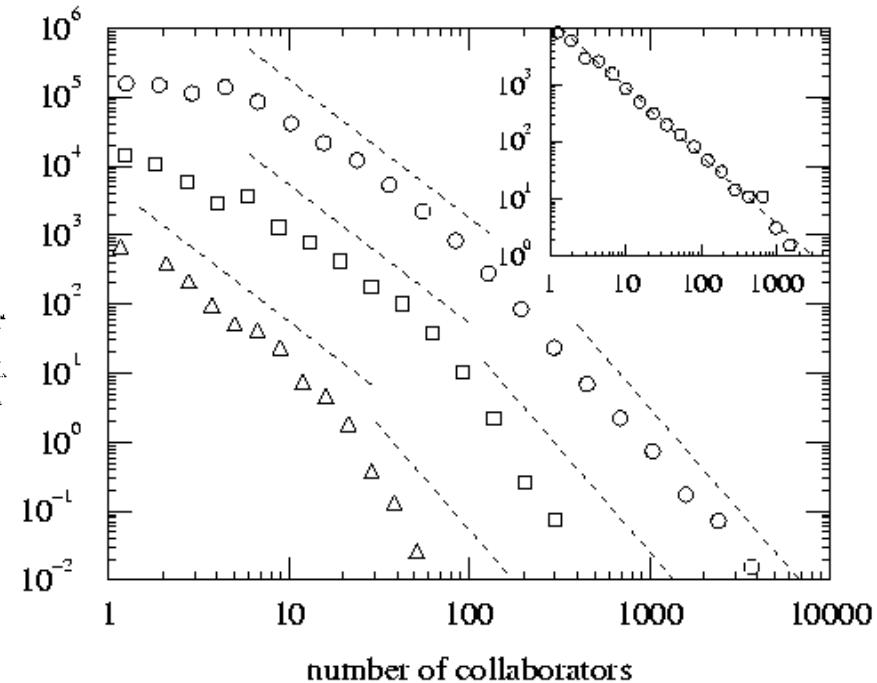
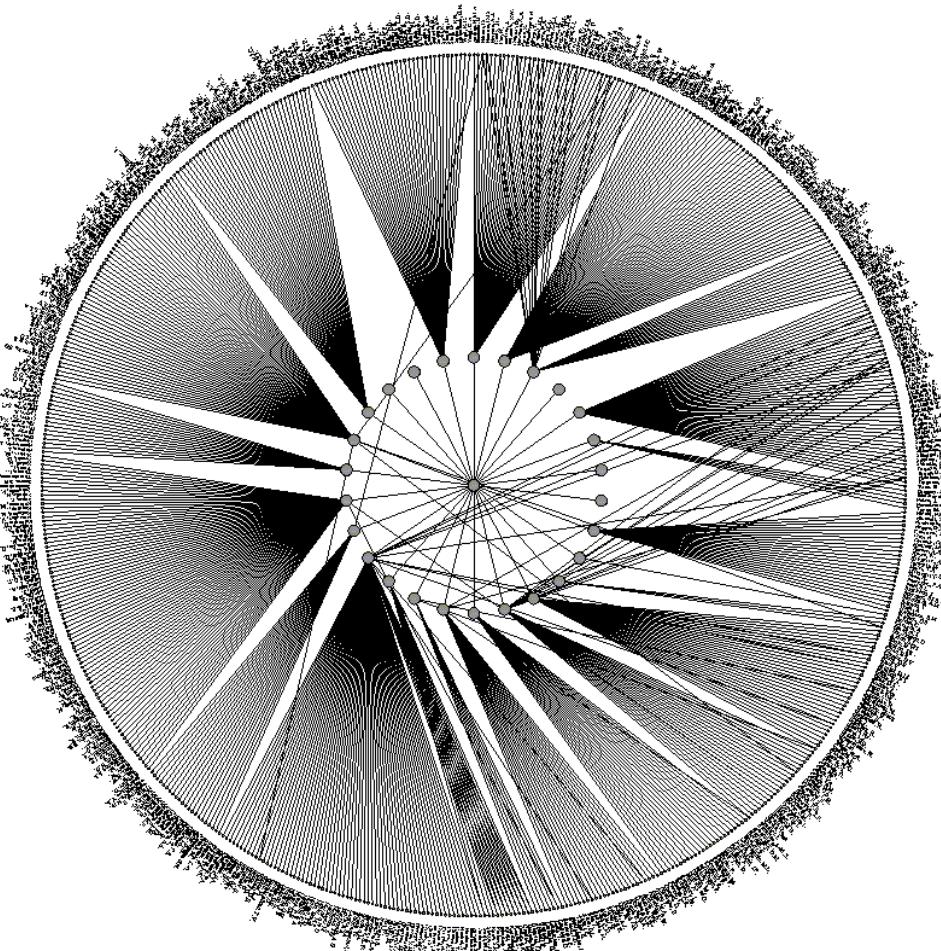
Scale-free Network



Science Coauthorship Network

Nodes: scientist (authors)

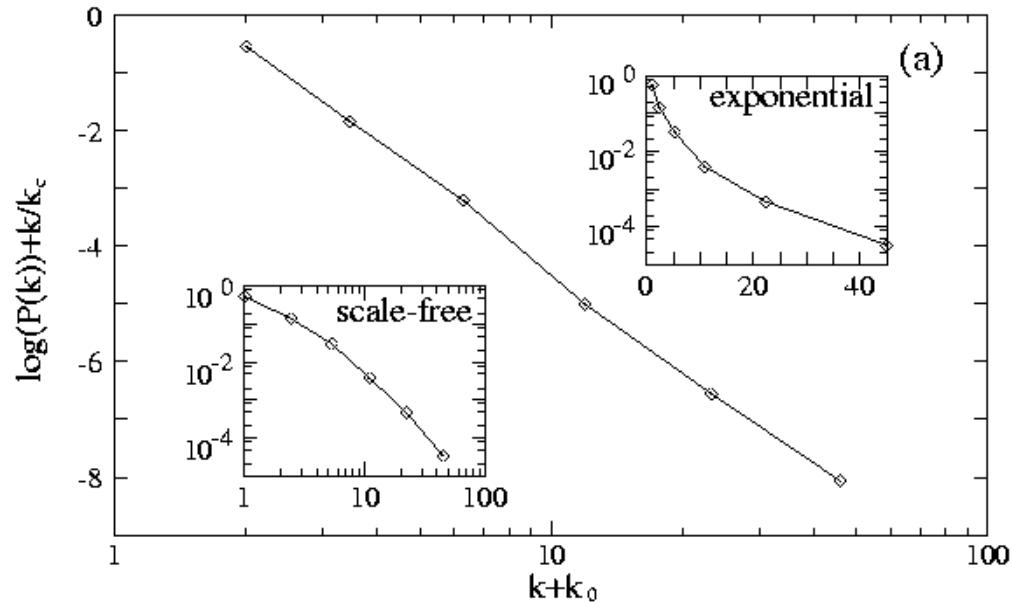
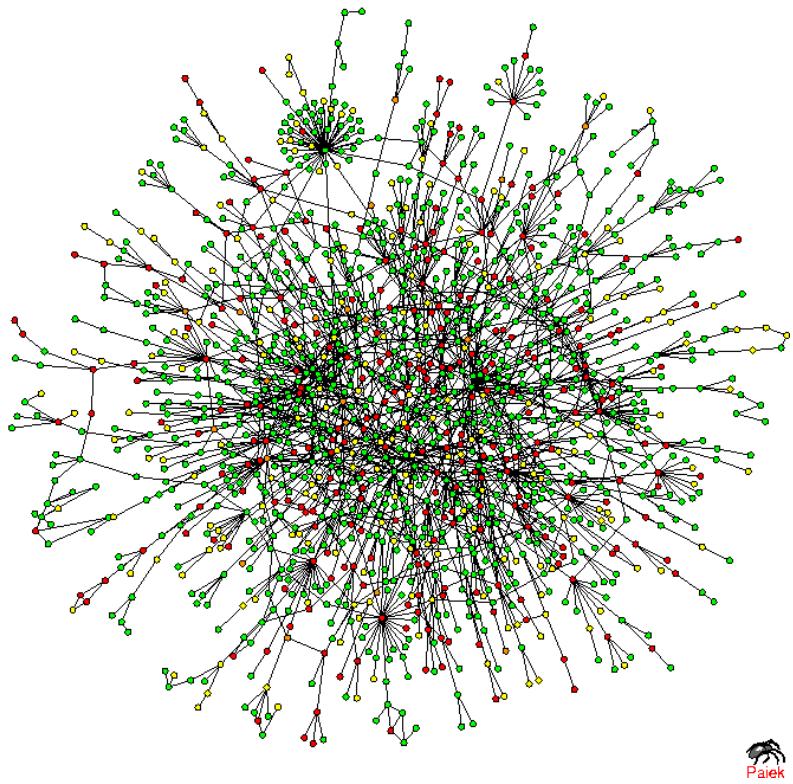
Links: write paper together



[Newman, 2000,

H. Jeong et al 2001]

Protein Protein Interaction network



$$P(k) \sim (k + k_0)^{-\gamma} \exp\left(-\frac{k + k_0}{k_\tau}\right)$$

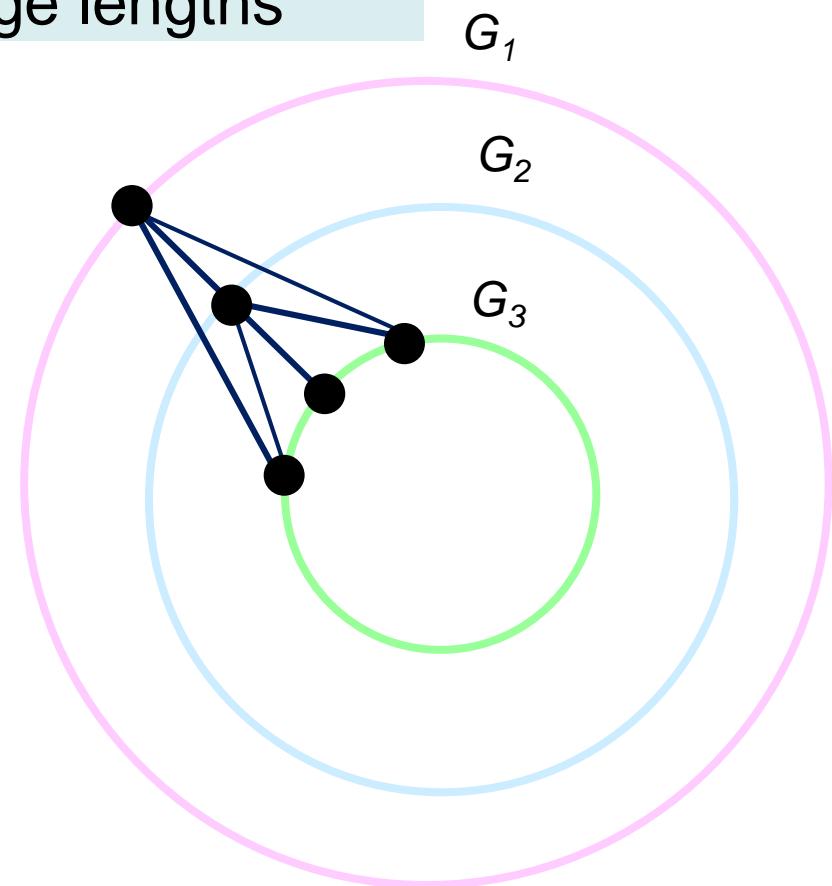
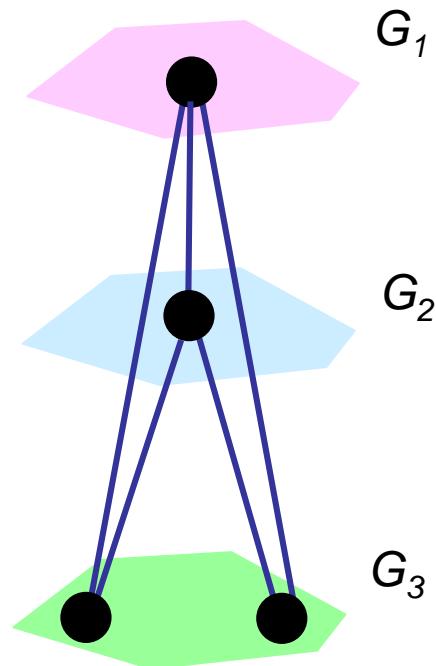
H. Jeong, S.P. Mason, A.-L. Barabasi, Z.N. Oltvai, Nature 411, 41-42 (2001)

Parallel Plane/Concentric Sphere Layout Criteria

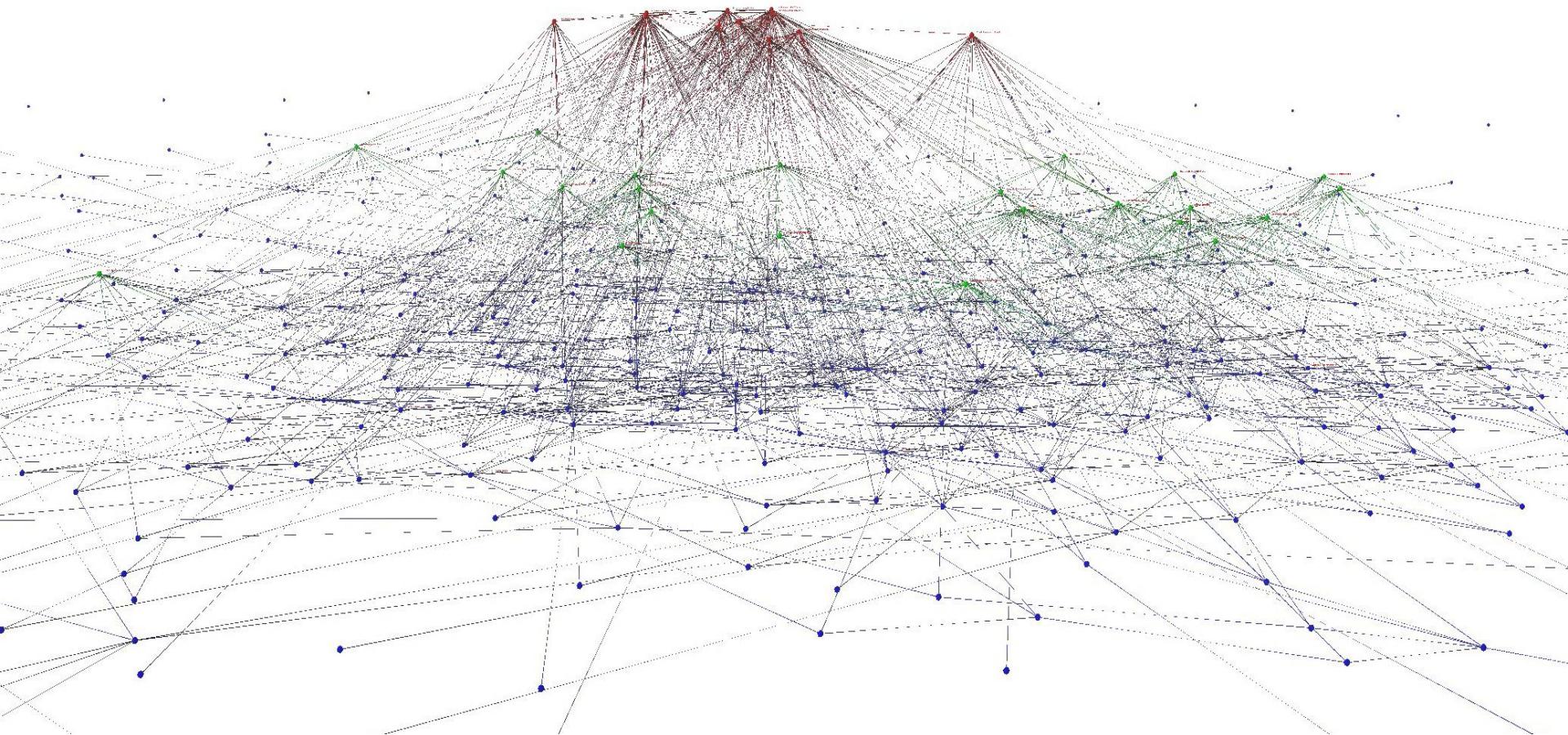
Criteria:

- Minimise occlusion
- Minimise intra-sphere edge crossings
- Minimise total inter-plane edge lengths

平行平面/同心球体布局标准
标准：
•尽量减少遮挡
•最小化球内的边交叉
•最小化总平面间边长



Plane layout: Citation Network

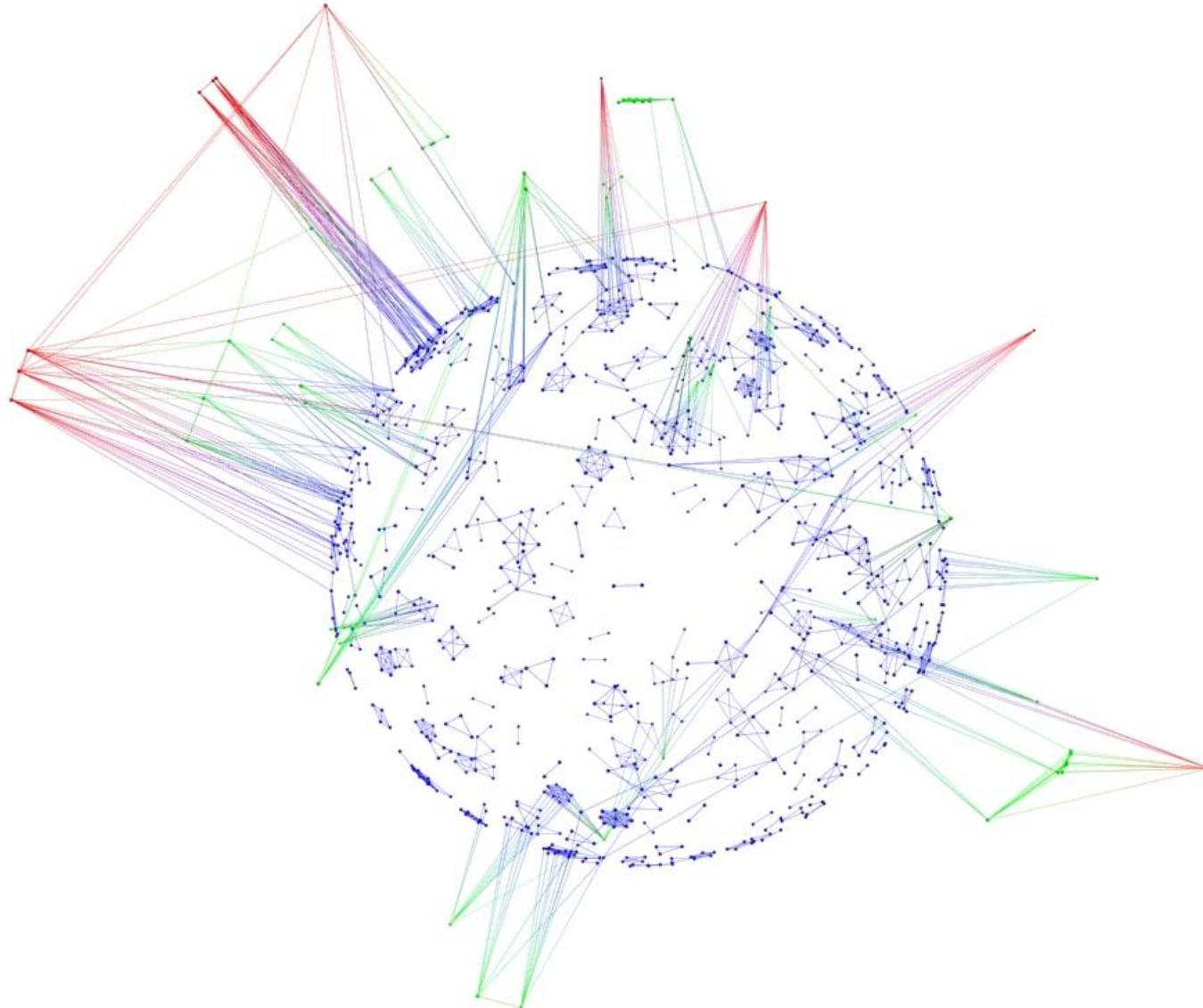


$|V|=535, |E|=1562 (>20, >10)$

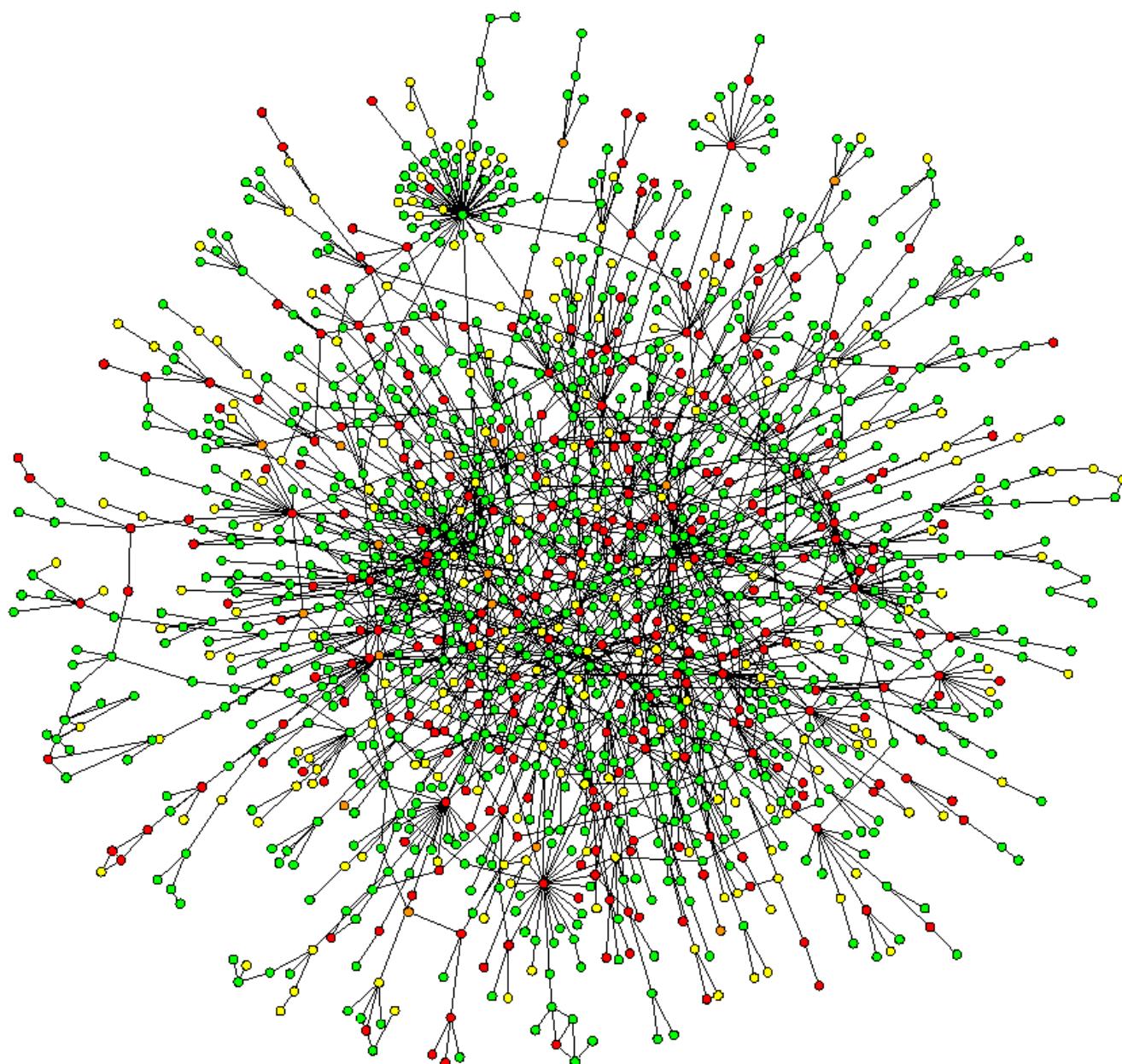
Apex: Highly Cited Papers



Sphere layout: Collaboration Network



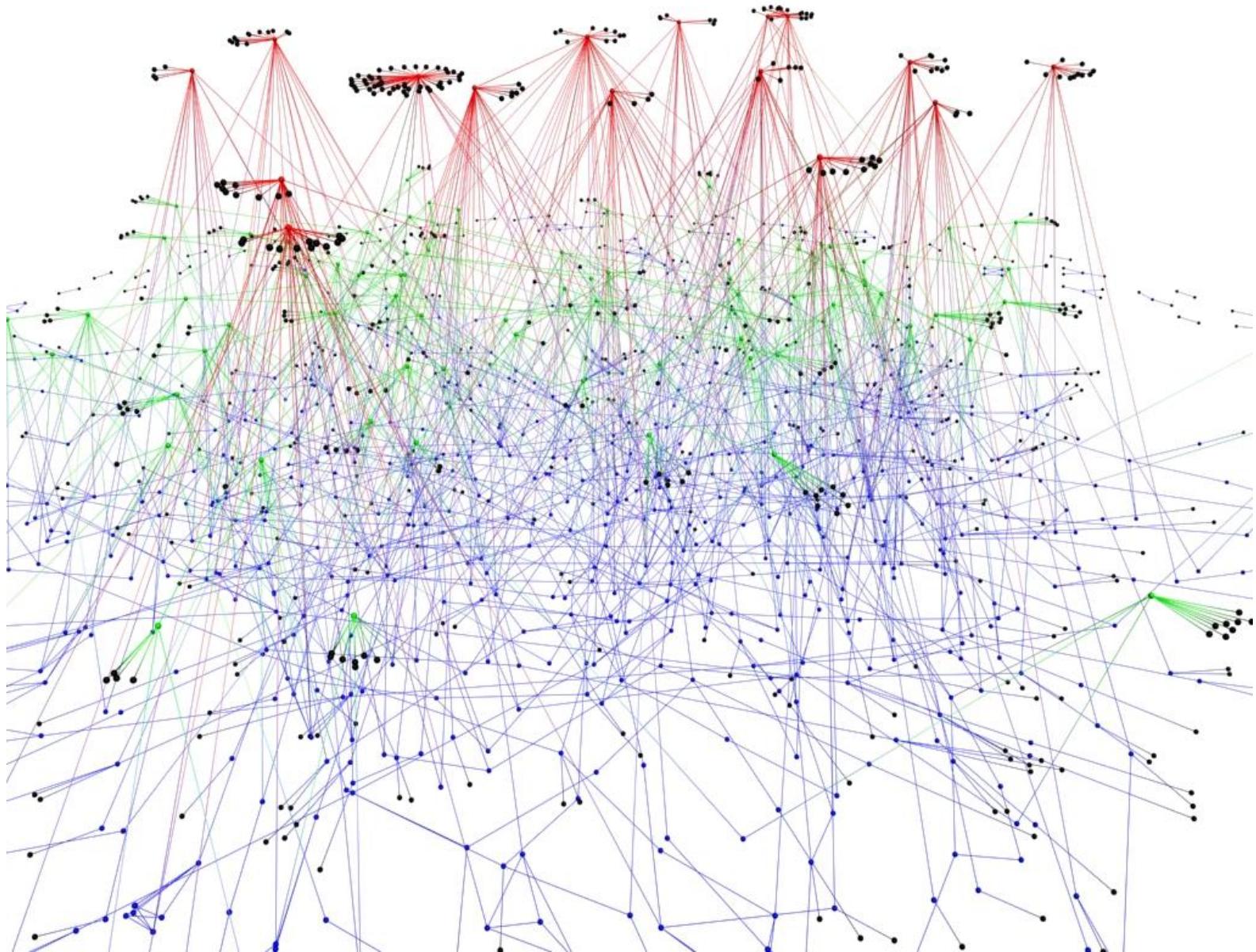
Apex: prominent researchers, $|V|=982$, $|E|=2012$
Research Empire: research clusters



PPI networks

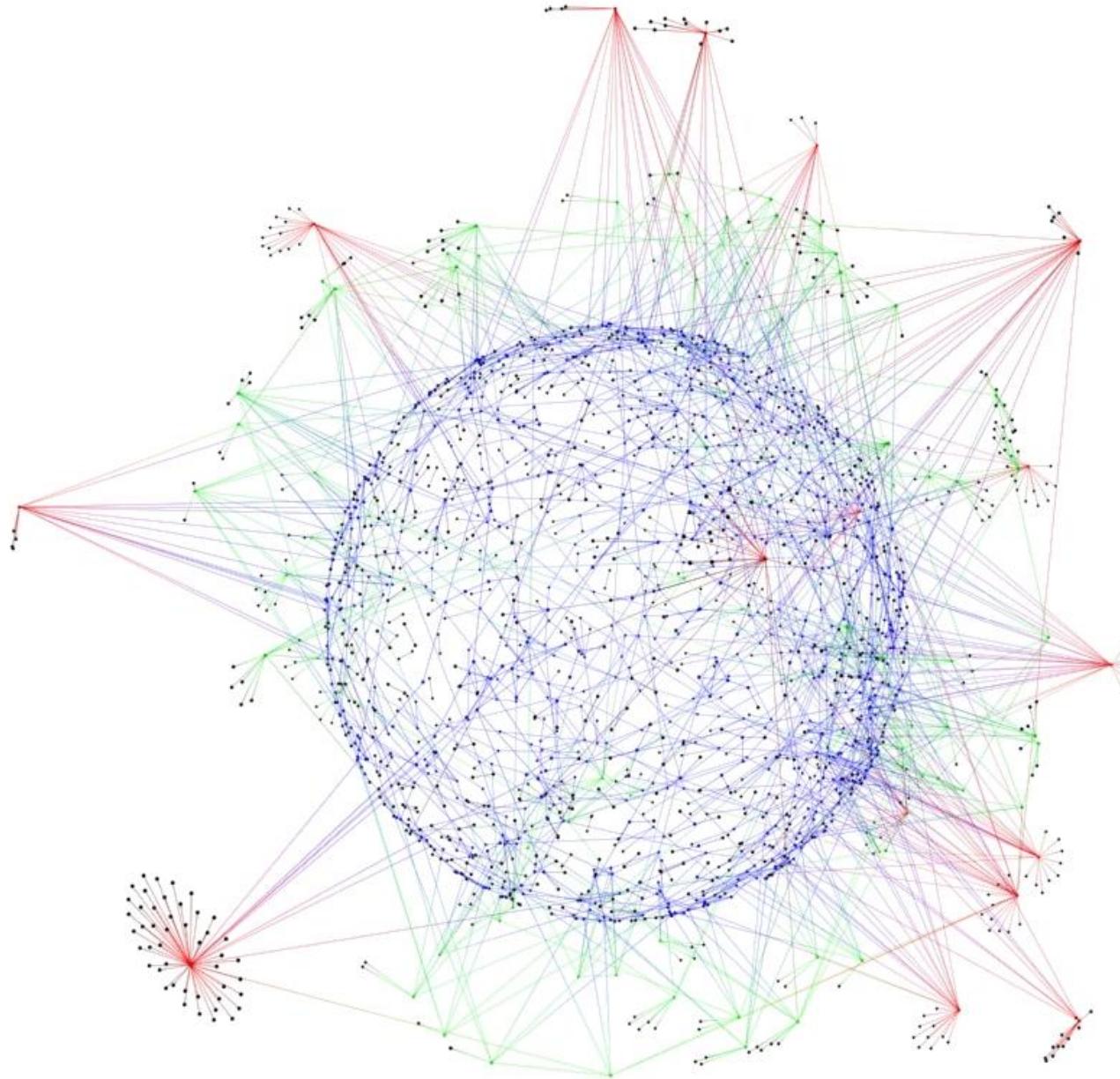
Hawoong Jeong

Plane layout: PPI Network



Apex: Lethal Proteins, $|V|=1846$, $|E|=2203$ ($>15, >7$)

Sphere layout: PPI Network

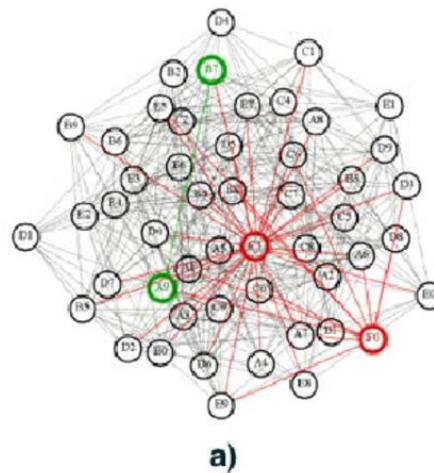


4. Reduce Visual Complexity

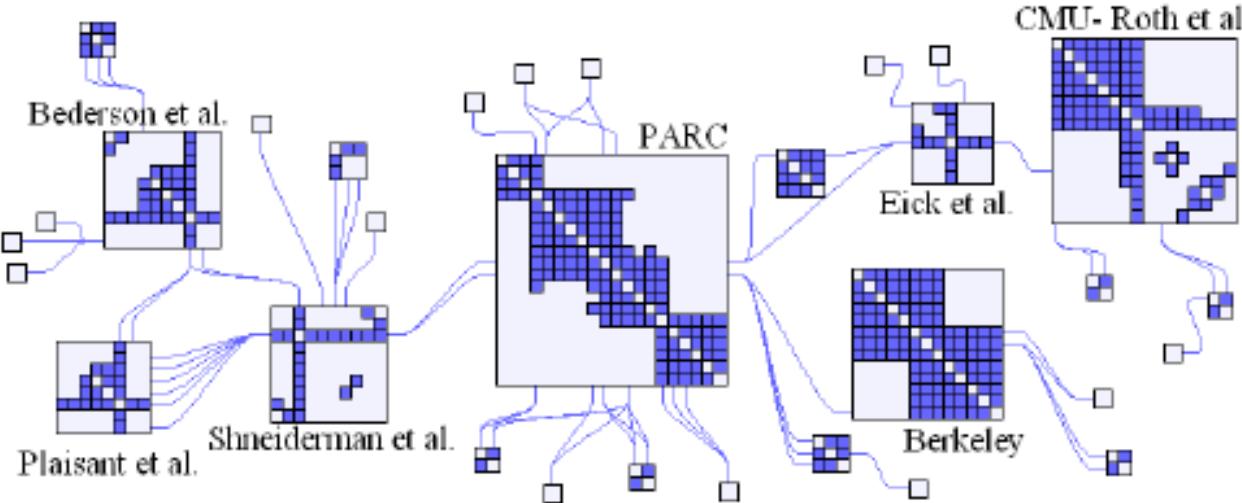
- *Matrix representation*
- *Map representation*
- *Edge bundling method*

Graph Representation

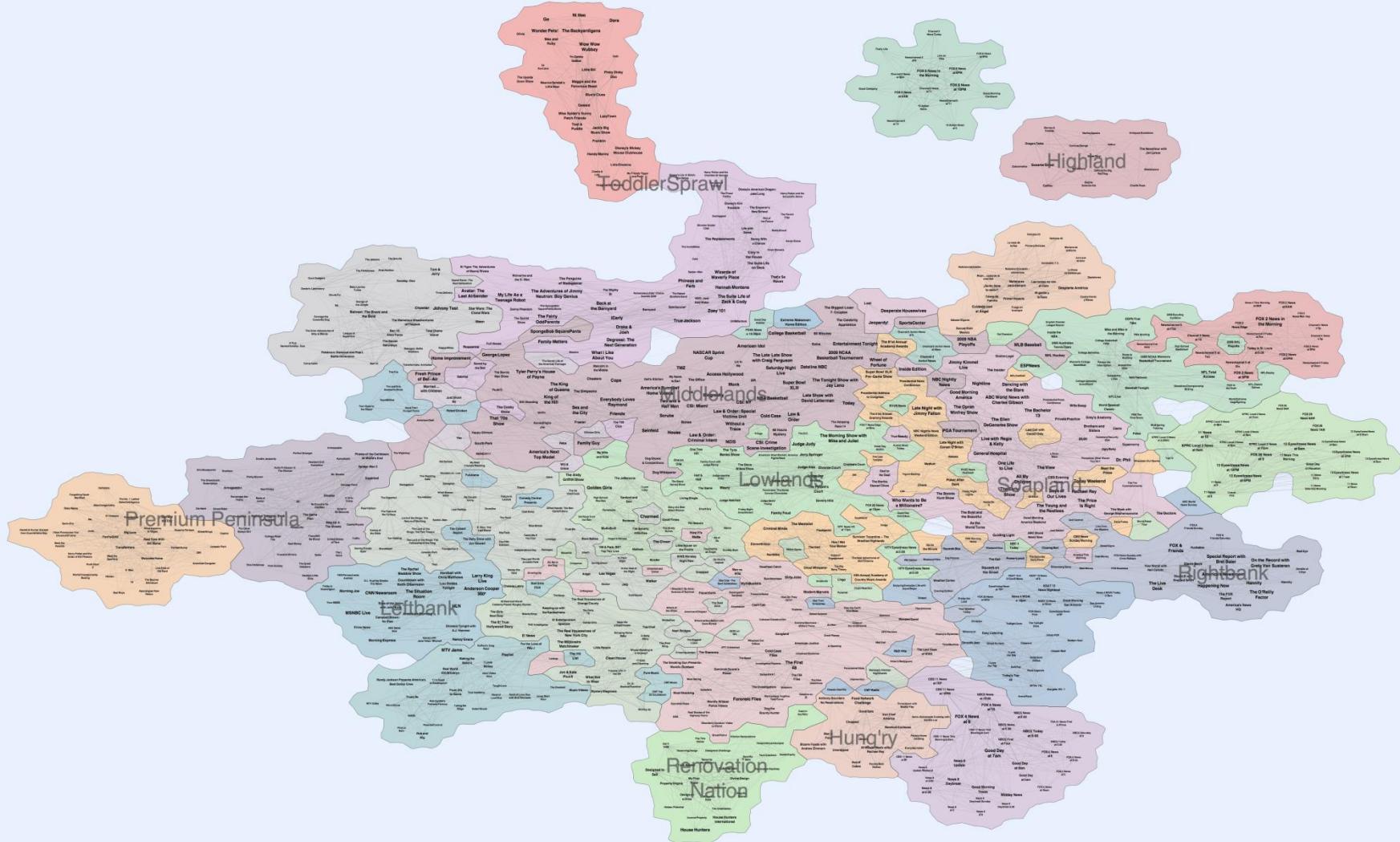
Node-link diagram



Matrix representation



Gmap [Ganser,Hu,Kobourov 2010]

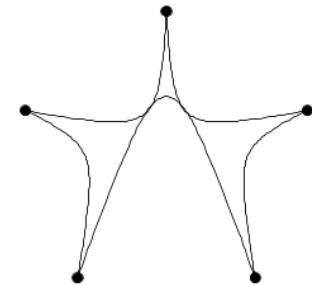
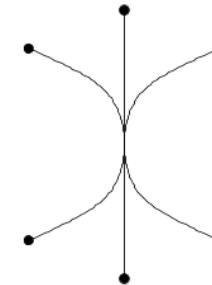


<http://yifanhu.net/MAPS/index.html>

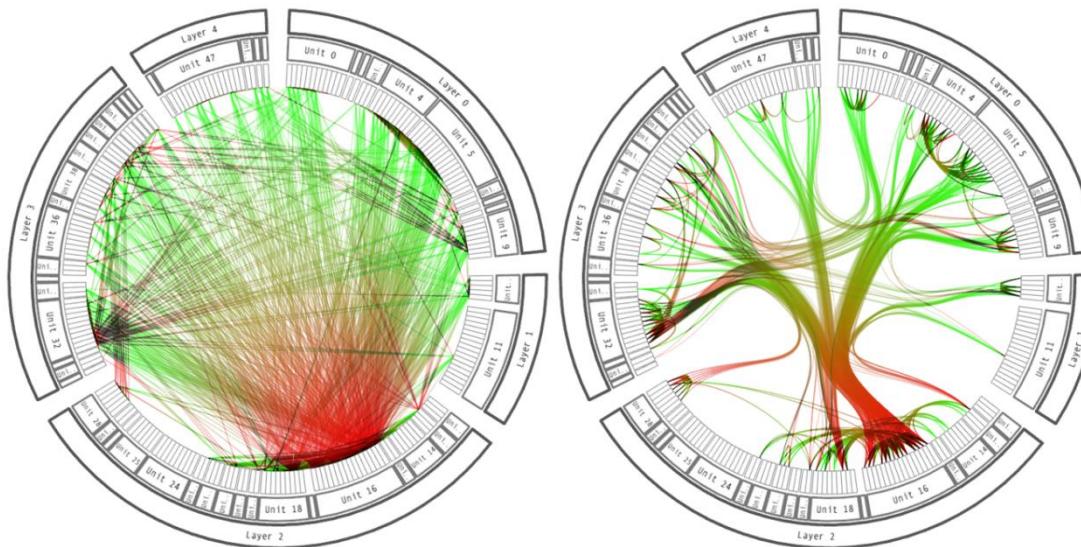
<https://www2.cs.arizona.edu/~kobourov/PROJECTS/maps.html>

Edge Bundling

- Edge Concentration [Newbery 1989]
- Confluent Drawing [Dickerson et al 2003]



- Edge Bundling [Holten 2006]



5. Integration with Analysis

- *Network Analysis*
- *Data Mining*
- *Machine Learning*

Social Network Analysis

- Methodological approach using graph theoretic concepts to describe, understand, explain social structure
- Purpose/level of interest
 - (1) **Centrality**: important actors/ crucial links
 - (2) **Cohesive subgroups**: components, cores, cliques
 - (3) **Structural roles**: positions, roles, clusters
 - (4) **Patterns**
 - (5) Network Statistics/Comparison

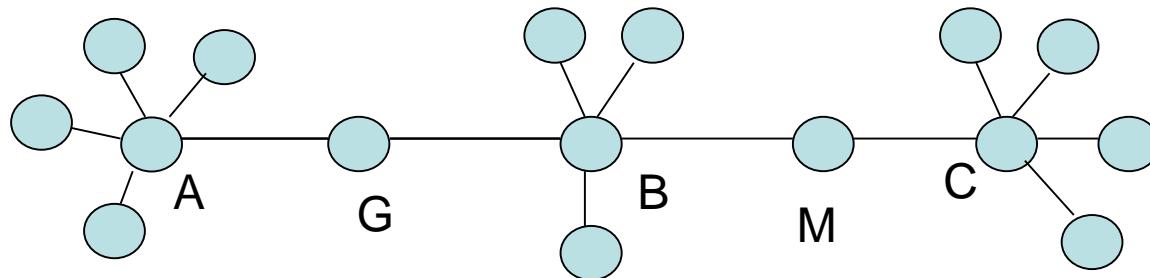
方法论方法使用图论理论来描述，理解，解释社会结构

• 目的/兴趣程度

- (1) 中心性：重要的参与者/关键环节
- (2) 有凝聚力的子群：组成部分，核心，团体
- (3) 结构角色：职位，角色，集群
- (4) 模式
- (5) 网络统计/比较

1. Centrality Analysis

- Degree: local measure
- Distance measures (global)
 - Betweenness [Freeman 79]
中间状态
 - low degree can be important (broker, gatekeeper, intermediary)
 - proportion of shortest paths connecting each pair X and Z which pass thru Y
 - Closeness: sum of shortest paths to all the other vertices
 - Eccentricity: length of the longest shortest path 离心率
- Feedback measures
 - Status/Hub/authority/eigenvector...

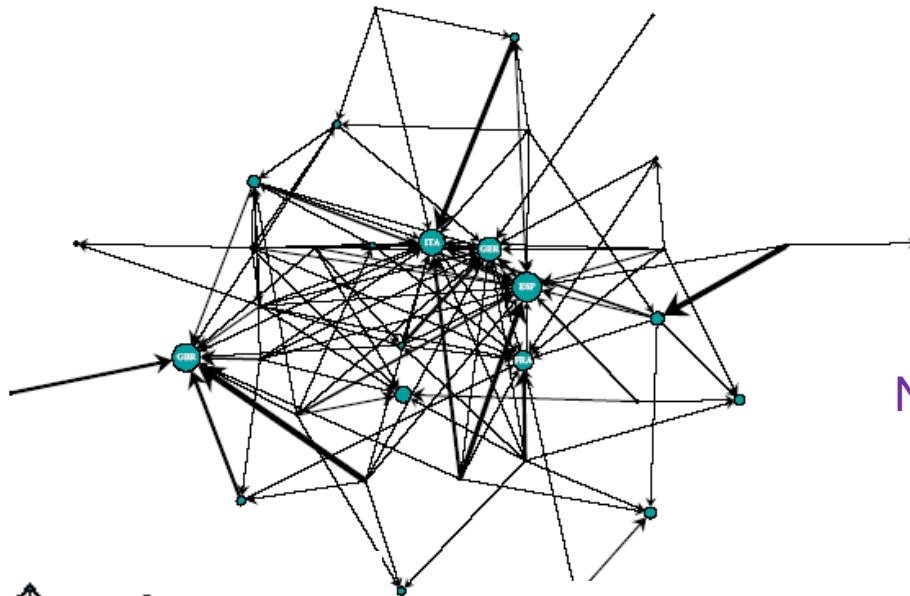


Degree: A, B, C
Betweenness: B, G, M

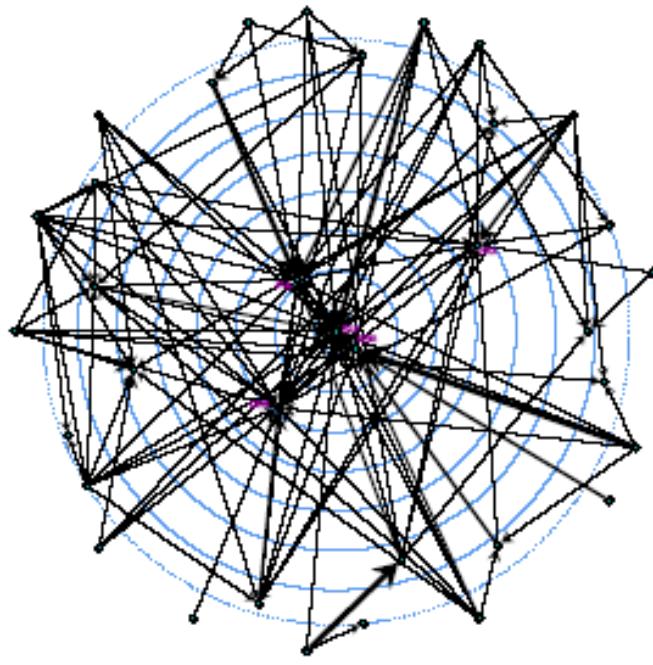
index	definition
local measures	
degree	$c_v = \sum_{e \in \text{instar}(v) \cup \text{outstar}(v)} \omega(e)$
indegree	$c_v = \sum_{e \in \text{instar}(v)} \omega(e)$
outdegree	$c_v = \sum_{e \in \text{outstar}(v)} \omega(e)$
distance measures	
betweenness	$c_v = \sum_{s \neq v \neq t \in V} \frac{\sigma_G(s, t v)}{\sigma_G(s, t)}$ where $\sigma_G(s, t)$ and $\sigma_G(s, t v)$ are the number of all shortest st -paths and those passing through v
closeness	$c_v = \frac{1}{\sum_{t \in V} \delta(v, t)}$
eccentricity	$c_v = \frac{1}{\max_{t \in V} \delta(v, t)}$
radiality	$c_v = \frac{\sum_{t \in V} (\text{diam}(G) + 1 - \delta(v, t))}{(n - 1) \cdot \text{diam}(G)}$

feedback measures	
status	$c_v = \alpha \cdot \sum_{(u,v) \in \text{instar}(v)} (1 + c_u)$ where $\alpha = \min\{\max_{v \in V} \text{indeg}(v), \max_{v \in V} \text{outdeg}(v)\}^{-1}$
eigenvector	$c_v = \mu^{-1} \sum_{(u,v) \in \text{instar}(v)} \omega(u, v) \cdot c_u$ where μ is the largest eigenvalue of $A(G)$
pagerank	$c_v = \gamma \cdot \frac{1}{n} + (1 - \gamma) \sum_{(u,v) \in \text{instar}(v)} c_u$ where $0 < \gamma < 1$ is a free parameter
authority	$c_v = \mu^{-1} \cdot \sum_{(u,v) \in \text{instar}(v)} \omega(u, v) \cdot \sum_{(u,w) \in \text{outstar}(u)} \omega(u, w) c_w$ where μ is the largest eigenvalue of $A(G)^T A(G)$
hub	$c_v = \mu^{-1} \cdot \sum_{(v,w) \in \text{outstar}(v)} \omega(v, w) \cdot \sum_{(u,w) \in \text{instar}(w)} \omega(u, w) c_u$ where μ is the largest eigenvalue of $A(G)A(G)^T$

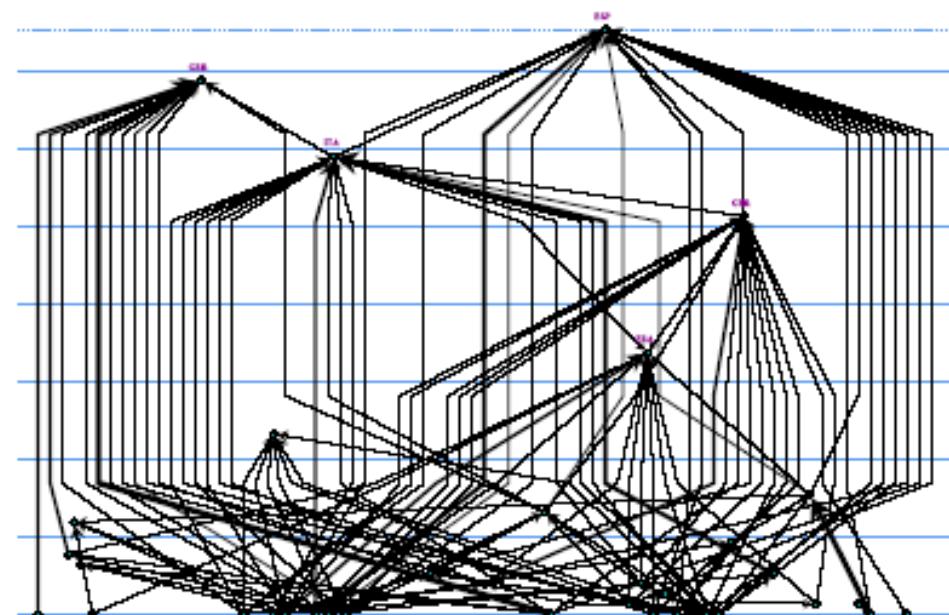
Centrality: Visualisation



Node size, edge weight



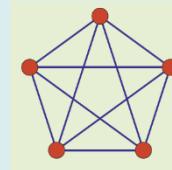
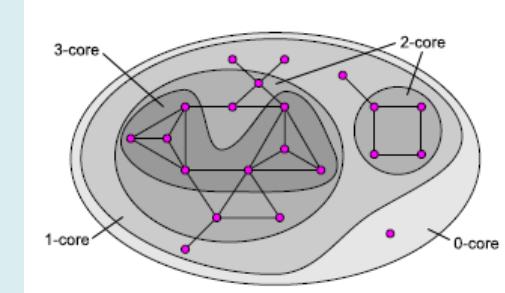
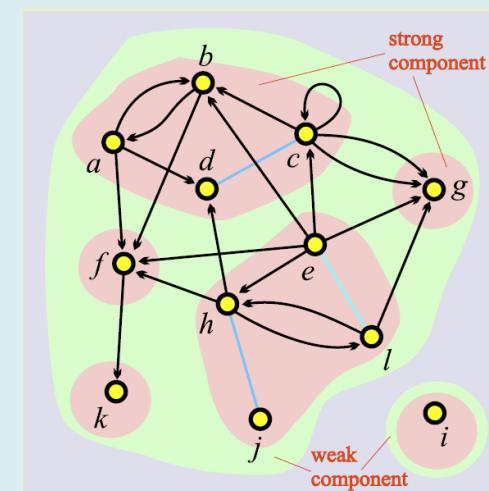
Radial drawing



Hierarchical drawing

2. Cohesive Subgroup

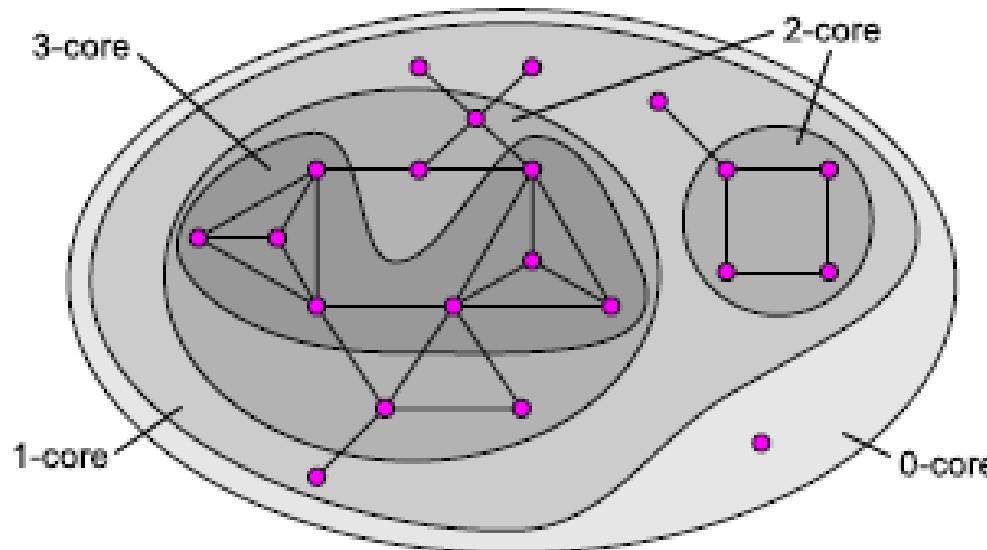
- **Component:**
 - Strong/weak component
 - Cycles/cyclic component
 - Connected/isolated (k -connectivity)
 - Cut vertex/ cut pair
- **k -Core**
 - maximal subgraph such that in which each vertex is adjacent to at least k other vertices
 - vertex degree: $\geq k$
- **Clique** : complete subgraph (strong/weak clique)
- **n -clique, n -clan, k -plex**



***k*-Core Analysis**

- A group analysis method from Social Network Analysis.
- Identify important dense subgraphs from big complex networks
- Linear-time algorithm [Batagelj].
- Degeneracy: in graph theory

社交网络分析中的群体分析方法。
•从大型复杂网络中识别重要的密集子图
•线性时间算法[Batagelj]。
•简并度:在图论中



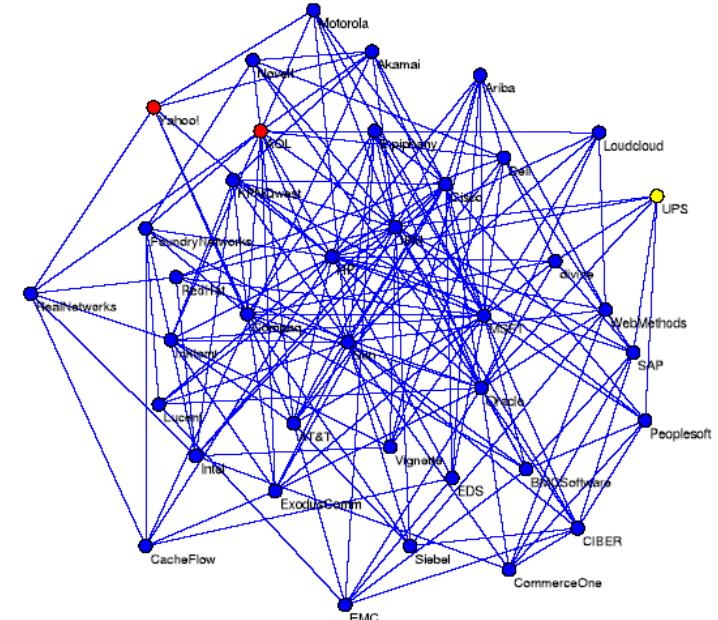
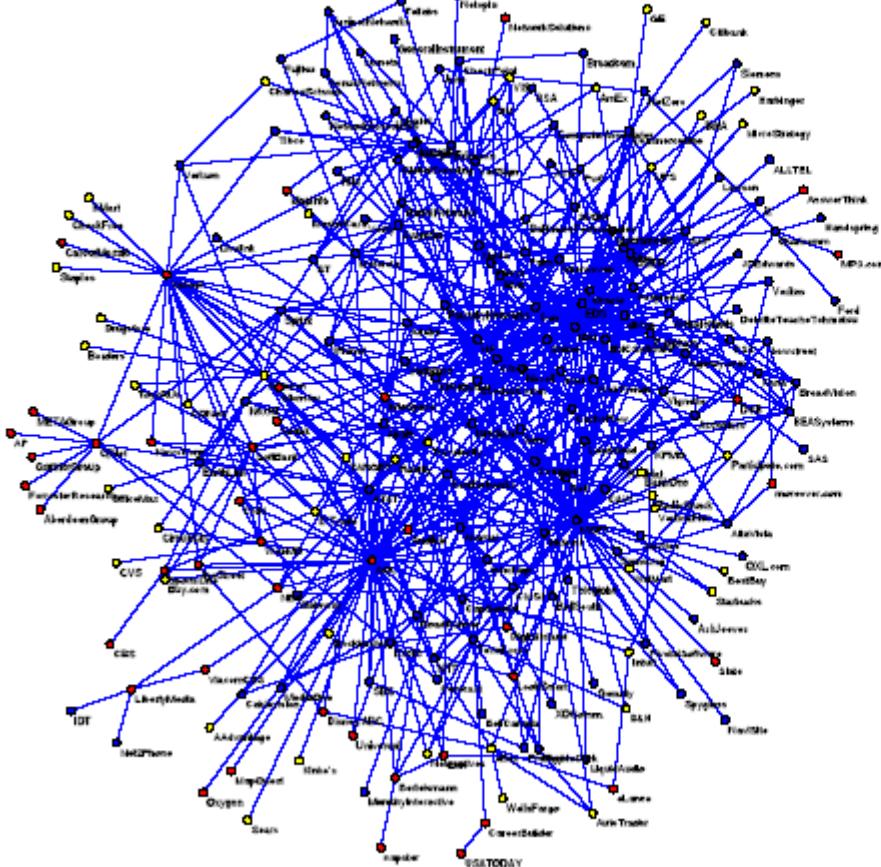
Example: K-core based Filtering

Krebs Internet industries

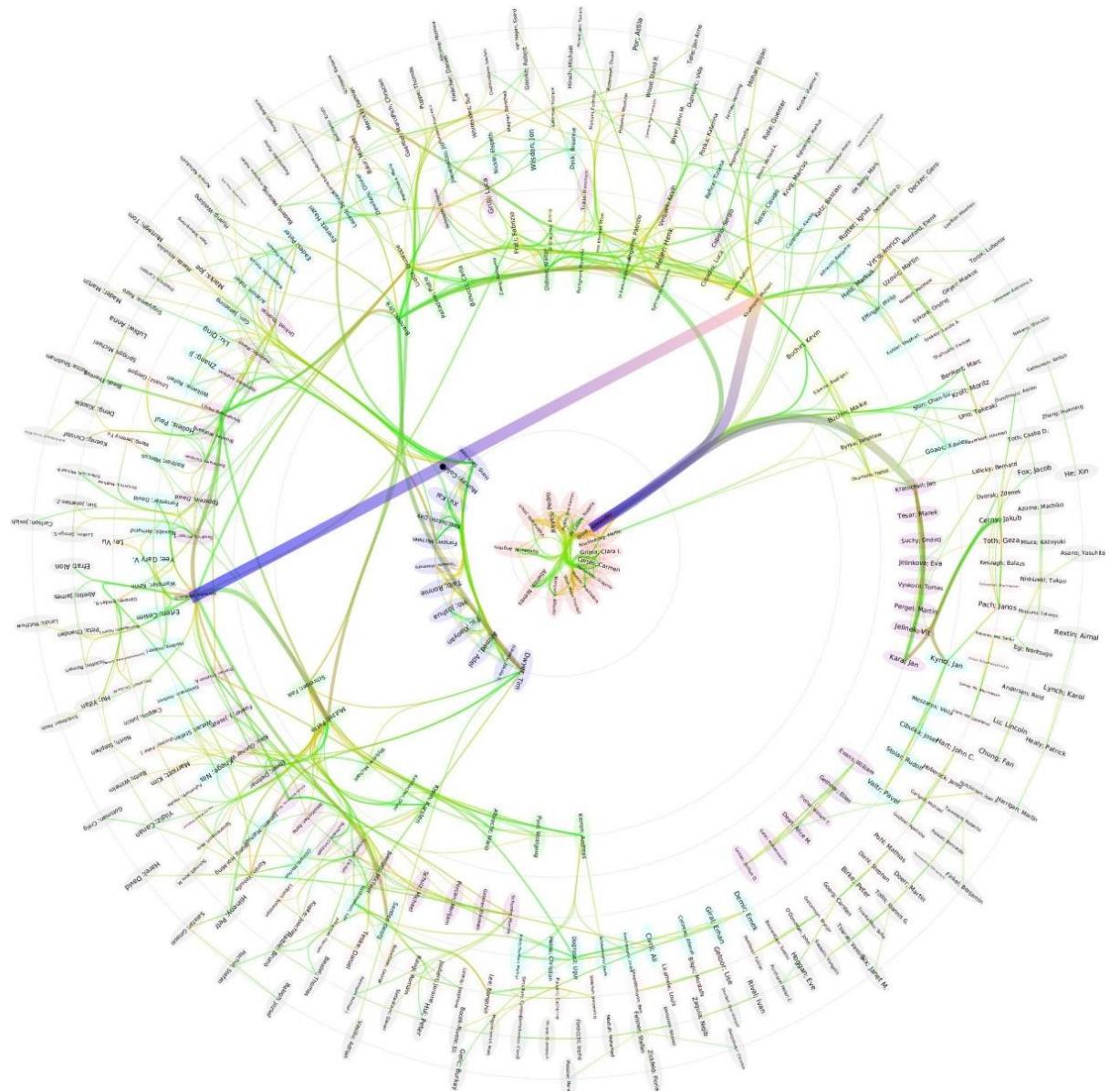
Each node in the network represents a company that competes in the Internet industry, 1998 do 2001.

$$n = 219, m = 631.$$

6-core of Krebs Internet industries



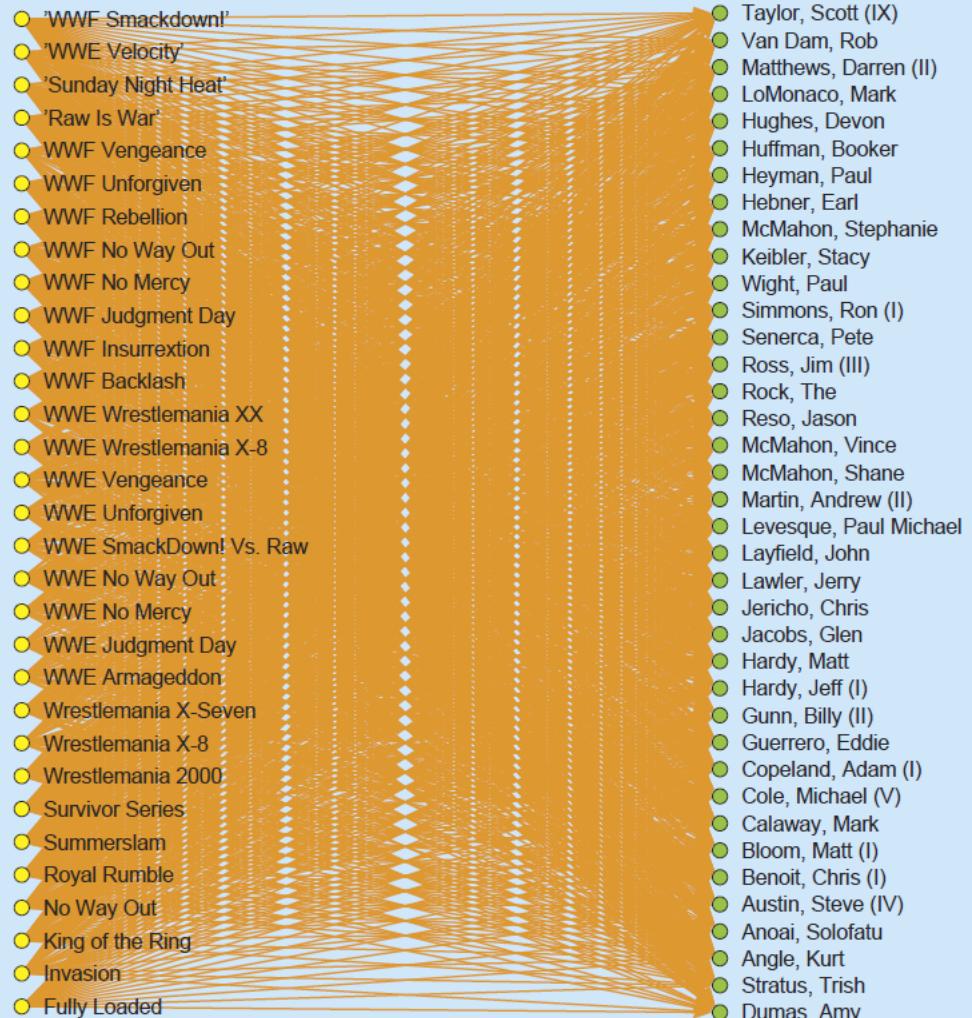
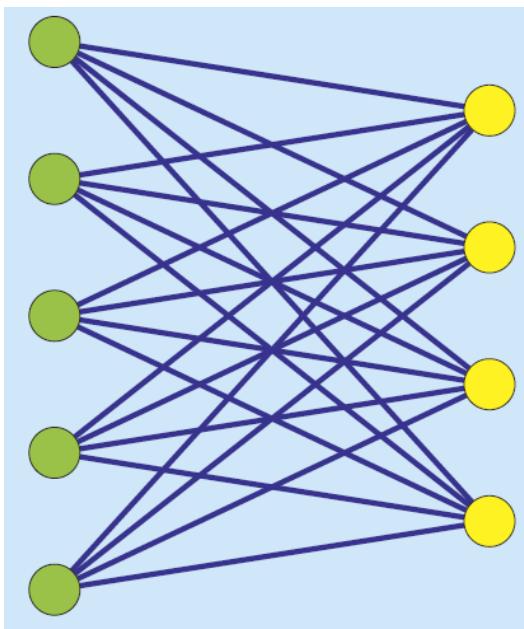
k-Core Visualisation: concentric circles



(p,q) -Core Analysis

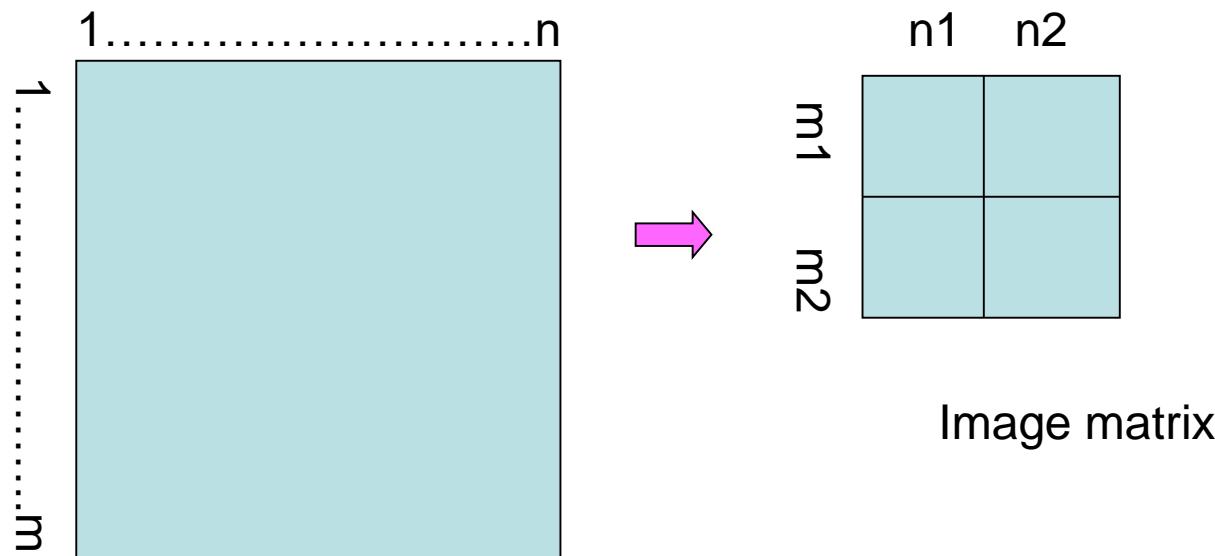
- Two-mode networks
- Bipartite graphs

二分图



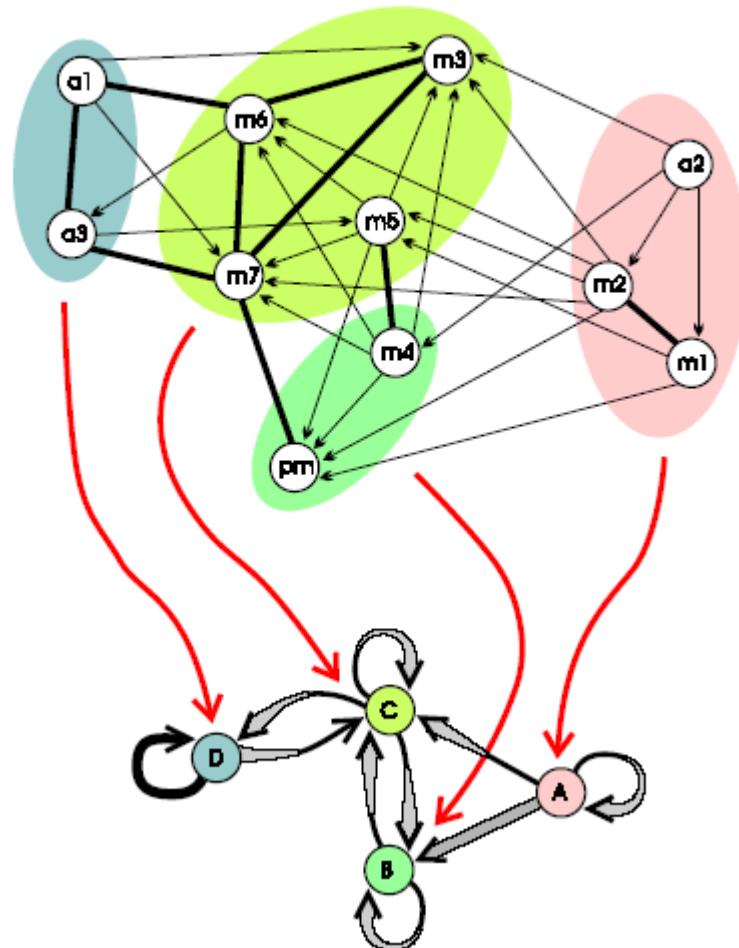
3. Network Positions and Structural Equivalence

- Block model (or image matrix) : reduction of complex network
- Structural/Automorphism/Regular Equivalence
- Clusters: cliques, distance, similarity



Blockmodeling as a clustering problem

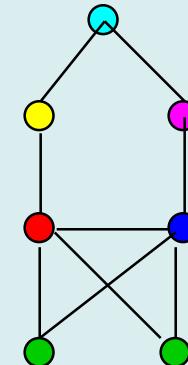
The goal of *blockmodeling* is to reduce a large, potentially incoherent network to a smaller comprehensible structure that can be interpreted more readily. Blockmodeling, as an empirical procedure, is based on the idea that units in a network can be grouped according to the extent to which they are equivalent, according to some *meaningful* definition of equivalence.



Structural Equivalence

The two green nodes are connected to exactly the same alters and are said to be *structurally equivalent*

They hold identical *positions* in the network



Formally, nodes a and b are *structurally equivalent* if, whenever (a,x) is an edge in G then so is (b,x) , and conversely ($x \neq a, b$)

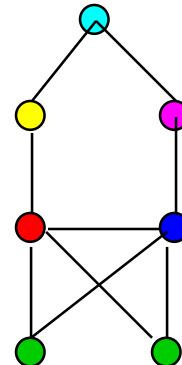
We can allocate nodes that are structurally equivalent to a (structural equivalence) *class*, or *block* – represented by colour in the diagram.

Blockmodel

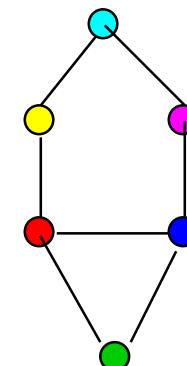
For a **structural equivalence** relation:

If a node in block x is connected to a node in block y , then **every** node in block x is connected to **every** node in block y

In this case we define an edge (x,y) between block x and block y , and the result is a **reduced graph** or **blockmodel** for the network



graph



blockmodel

Automorphic Equivalence

A **permutation** α of the node set V of a graph is a re-labelling of the nodes:

we write $\alpha(a)$ for the node to which the label a is attached by the relabelling

图的节点集V的置换 α 是节点的重新标记：

我们为标签所在的节点写 $\alpha(a)$

a 通过重新标记附加

A permutation α is an **automorphism** of the graph G if, whenever (a,x) is an edge in G , then $(\alpha(a),\alpha(x))$ is an edge in G , and conversely.

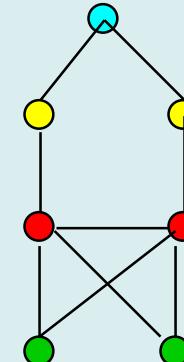
置换 α 是图G的自同构，如果每当 (a, x) 是G中的边，则 $(\alpha(a), \alpha(x))$ 是G中的边，并且相反。

Nodes a and b are **automorphically equivalent** if $b = \alpha(a)$ for some automorphism α

如果对于某些自同构 α , $b = \alpha(a)$ ，则节点a和b是自同等效的

示例：图中具有相同颜色的节点是自动等效的

Example: nodes with the same colour in the graph are automorphically equivalent

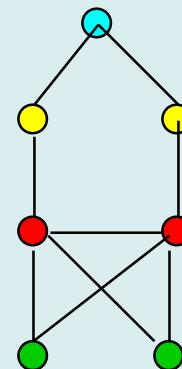


Why is automorphic equivalence interesting?

Automorphically equivalent nodes have the same position in a network in a more abstract sense than structurally equivalent nodes: they are not connected to the exact same nodes, but to nodes that play analogous roles in the network

a potential representation for roles: leader, principal, broker etc

Note: if nodes are structurally equivalent, then they are also automorphically equivalent, but the converse does not hold



graph



blockmodel

Regular Equivalence

Two nodes a and b are regularly equivalent if

- (a) whenever (a, x) is an edge in G , then there is some node y that is regularly equivalent to x for which (b, y) is an edge in G .

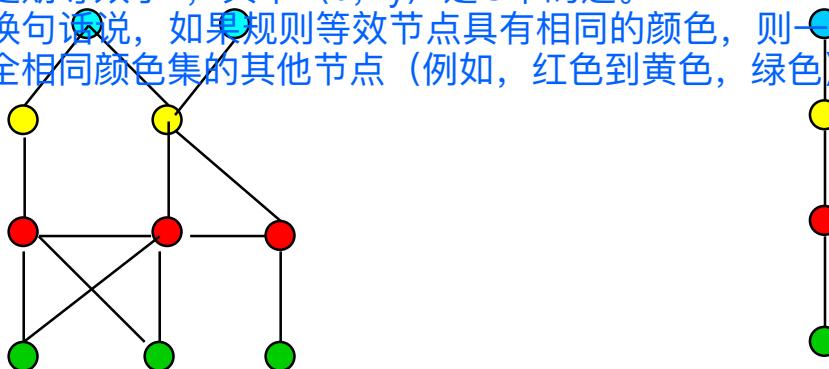
In other words, if regularly equivalent nodes have the same colour, then each node in a class of regularly equivalent nodes is connected to other nodes of exactly the same set of colours (e.g red to yellow, green)

如果，两个节点a和b通常是等效的

(a) 只要 (a, x) 是 G 中的边，那么就有一些节点 y

定期等效于 x ，其中 (b, y) 是 G 中的边。

换句话说，如果规则等效节点具有相同的颜色，则一类规则等效节点中的每个节点连接到完全相同颜色集的其他节点（例如，红色到黄色，绿色）

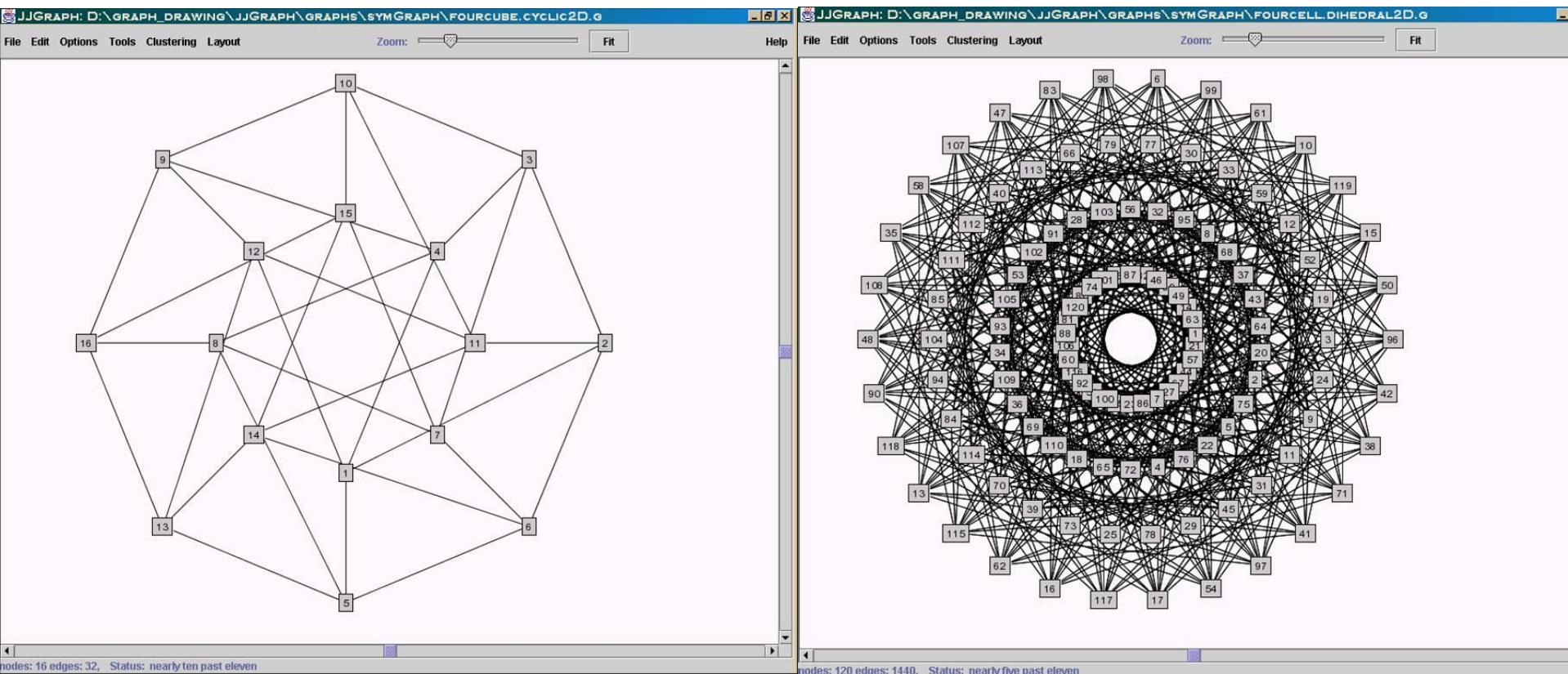


graph

(regular equivalence) blockmodel

Automorphic Equivalence: Visualisation

Symmetric drawing [Abelson,Hong,Taylor 03]

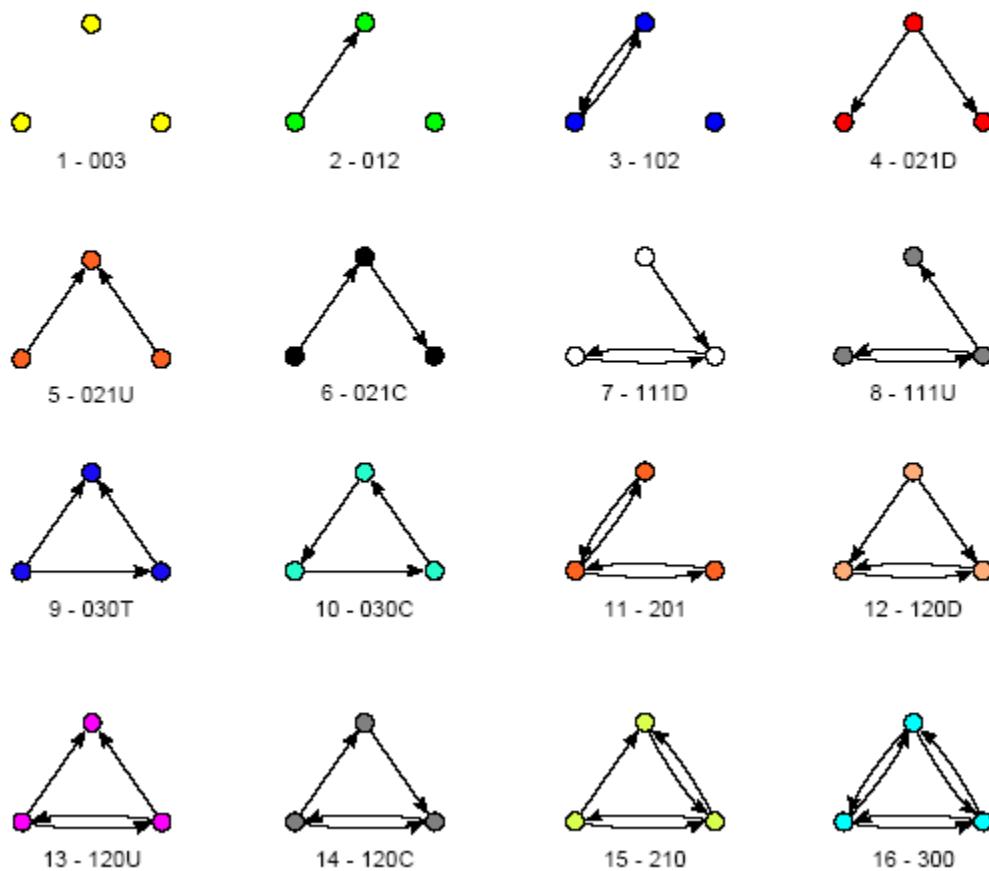


Concentric circle placement: based on orbit (equivalent class)

4. Patterns: Triads, Diads, Motifs

三联, 二连, 图案

Triads



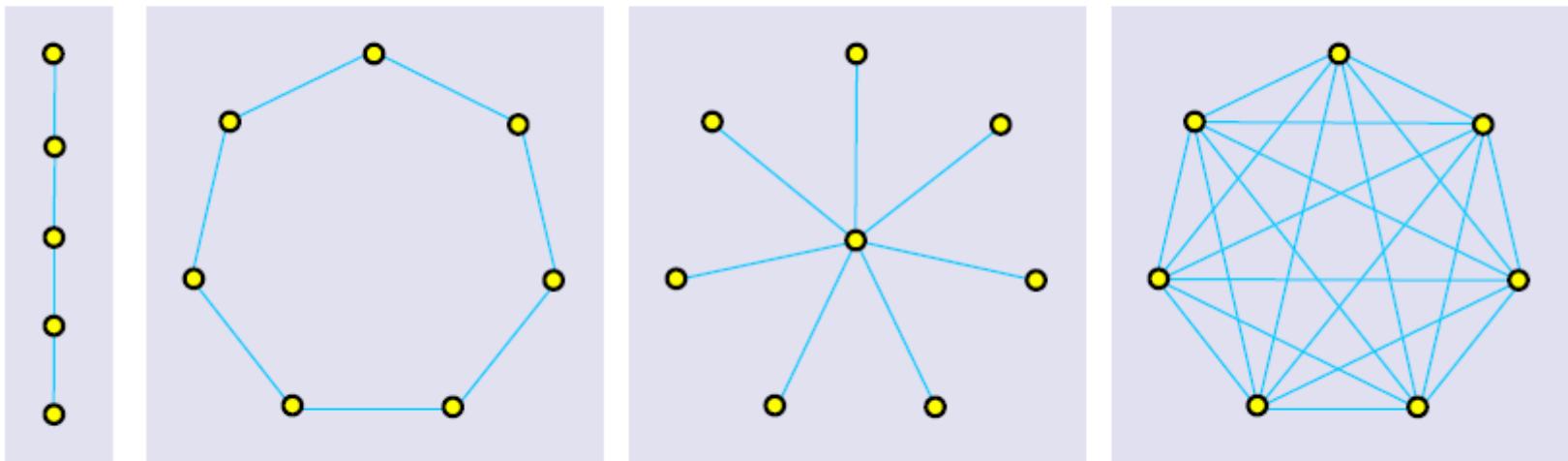
Let $\mathbf{G} = (V, R)$ be a simple directed graph without loops. A *triad* is a subgraph induced by a given set of three vertices.

There are 16 nonisomorphic (types of) triads. They can be partitioned into three basic types:

- the *null* triad 003;
- *dyadic* triads 012 and 102; and
- *connected* triads: 111D, 201, 210, 300, 021D, 111U, 120D, 021U, 030T, 120U, 021C, 030C and 120C.

Subgraph Patterns

Special graphs – path, cycle, star, complete



Graphs: *path* P_5 , *cycle* C_7 , *star* S_8 in *complete graph* K_7 .

5. Network Measures/Statistics

- **Degree distribution**
- **Clustering coefficient**
- **Diameter**
- **Average path length**
- **Connected component**
- **Density**

References

- **Social Network Analysis: Methods and Applications**
[Wasserman and Faust 94]
- **Network Analysis: Methodological Foundations**
LNCS 3418 Tutorial [Brandes and Erlebach eds. 04]
<http://dblp.uni-trier.de/db/conf/dagstuhl/na2004.html> (book chapters)

Social Network Analysis Tools

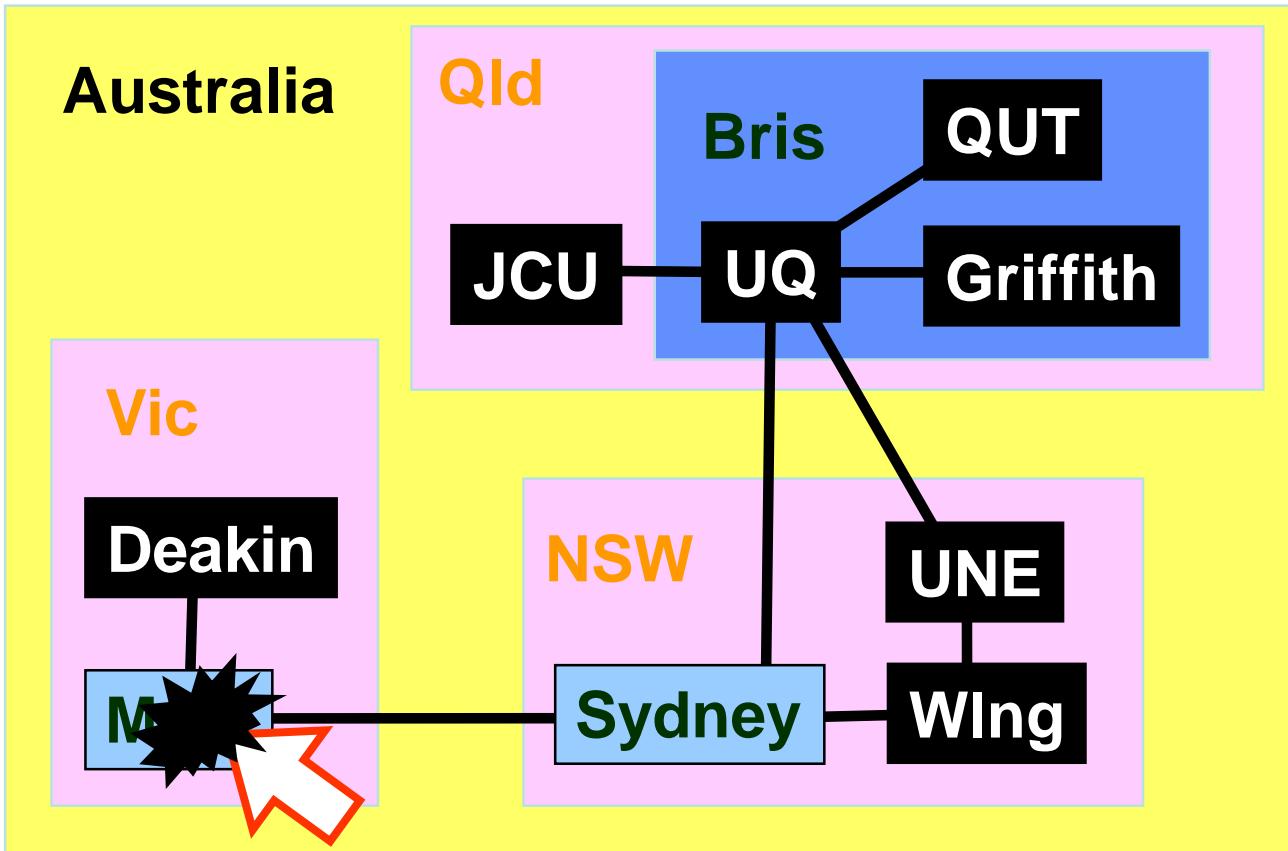
- **Pajek:** <http://mrvar.fdv.uni-lj.si/pajek/>, <http://mrvar.fdv.uni-lj.si/pajek/>
<http://vlado.fmf.uni-lj.si/pub/networks/doc/sunbelt/sunbeltxviii.pdf>
- **Visone:** <https://visone.info/>
- **NetMiner:** <http://www.netminer.com/main/main-read.do> (tutorial)

6. Interaction and Navigation

Interaction

- Ben Schneiderman mantra:
“Overview first, Zoom-in; Details on Demand”
- Types of interaction
 - Filtering
 - Selection
 - Zooming/panning
 - Collapsing/expanding
- Taxanomy [Yi et al. 2007]
 - **Select**: mark something as interesting
 - **Explore**: show me something interesting
 - **Reconfigure**: show me a different arrangement
 - **Encode**: show me a different representation
 - **Abstract/Elaborate**: show me more or less detail
 - **Filter**: show me something conditionally
 - **Connect**: show me related items

Clustered Graph [Pulo, Eades 03]



Australia

Vic

Deakin

Melb

LaTrobe

MU

Monash

Qld

JCU

Bris

QUT

UQ

Griffith

NSW

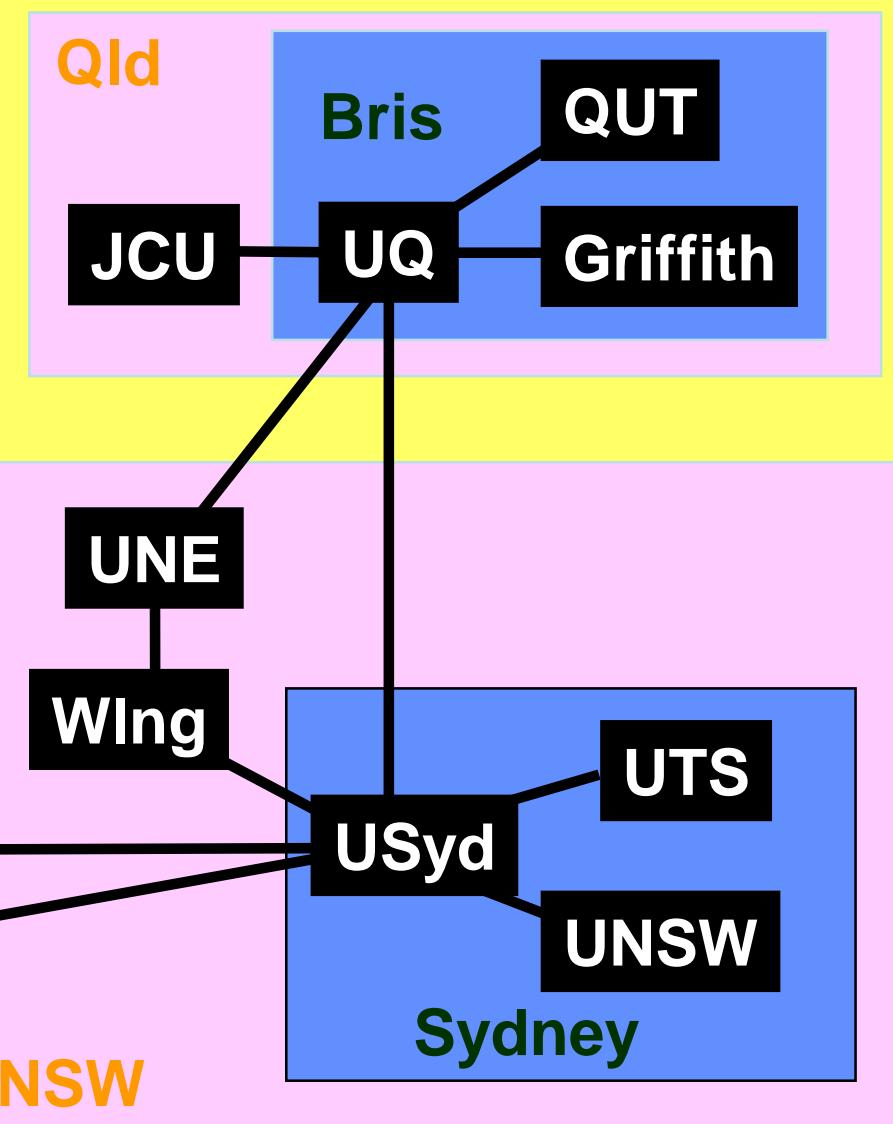
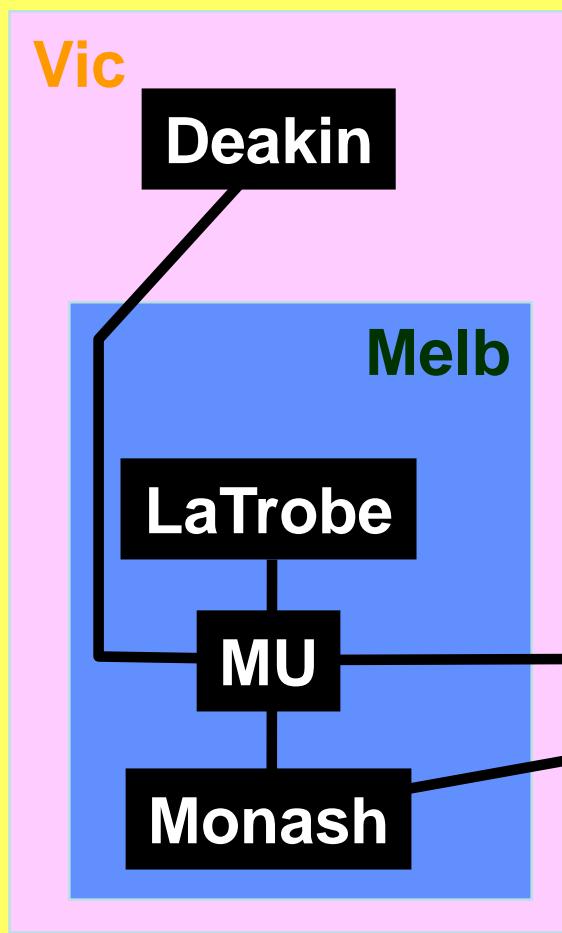
Sydn

UNE

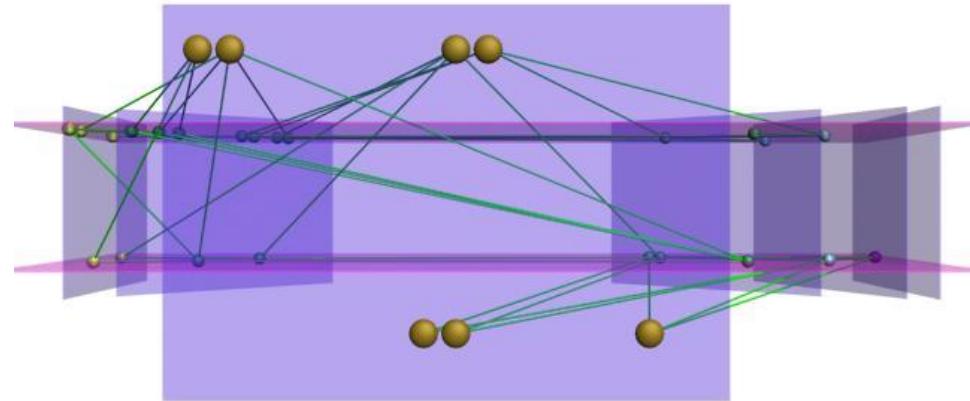
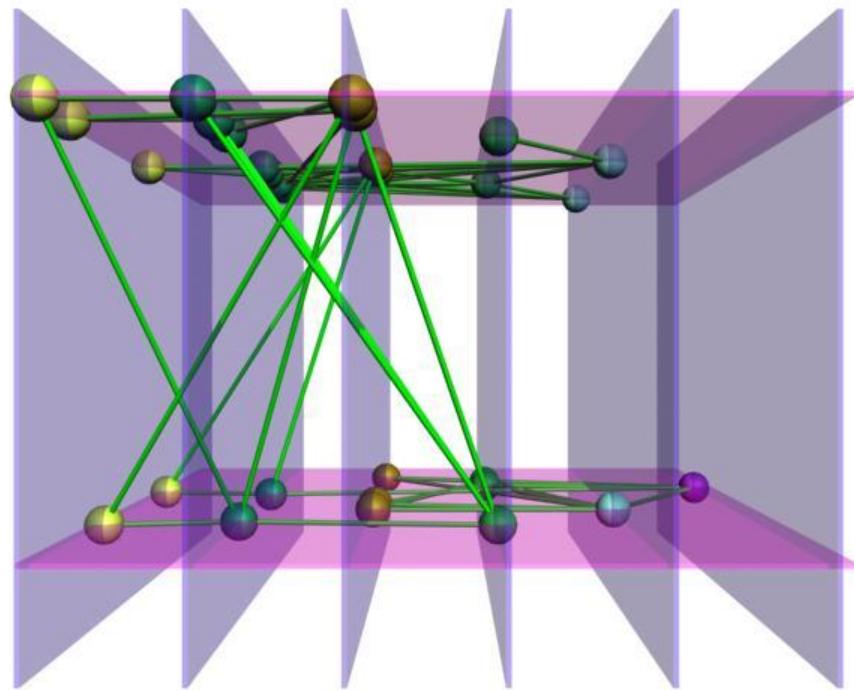
Wlng



Australia



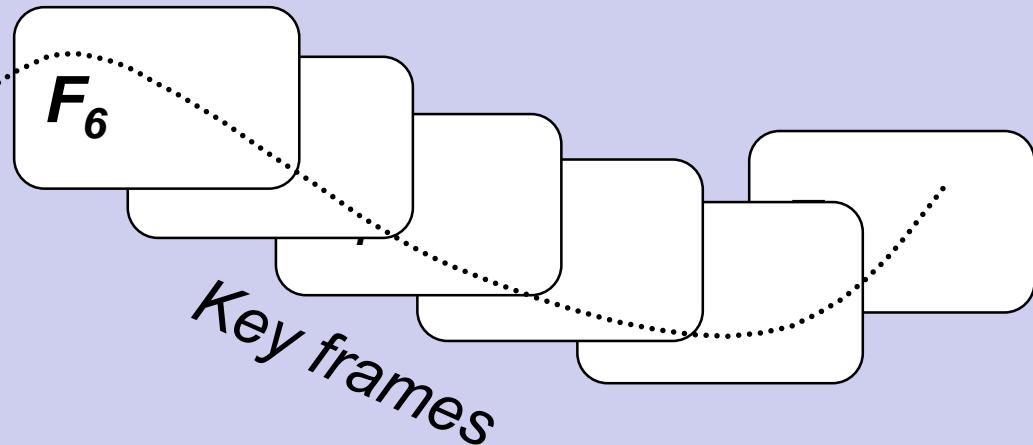
2.5D Graph Navigation [Ahmed, Hong 07]



Navigating Huge Unknown Network

Huge network

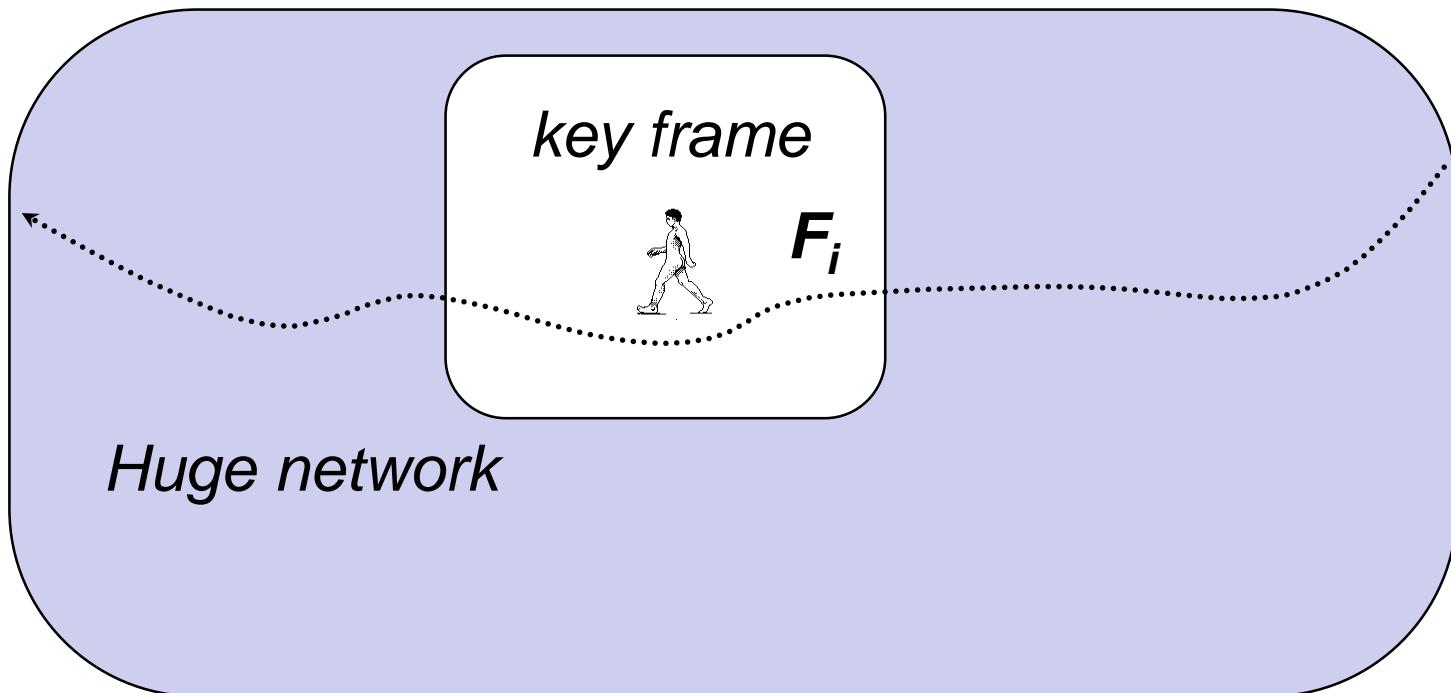
Path of exploration



Navigating Huge Unknown Network

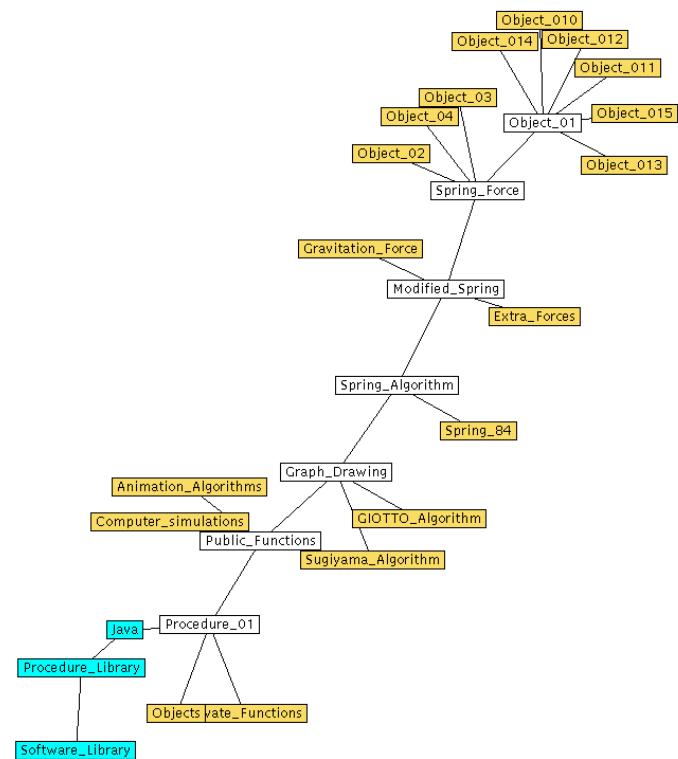
- We visualize a tiny part (a “key frame” F_i) of a huge network at time i .
- We move from F_i to F_{i+1} by user interaction.
- F_i overlaps F_{i+1} .

• 我们在第*i*时刻看到一个巨大网络的一小部分(“关键帧” F_i)。
• 我们通过用户交互从 F_i 移动到 F_{i+1} 。
• F_i 重叠 F_{i+1} 。



DA-TU (Maolin Huang [Hunag, Eades 1998])

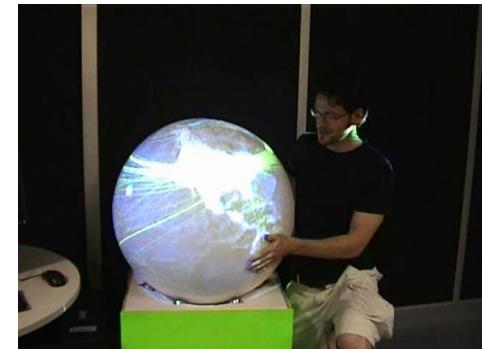
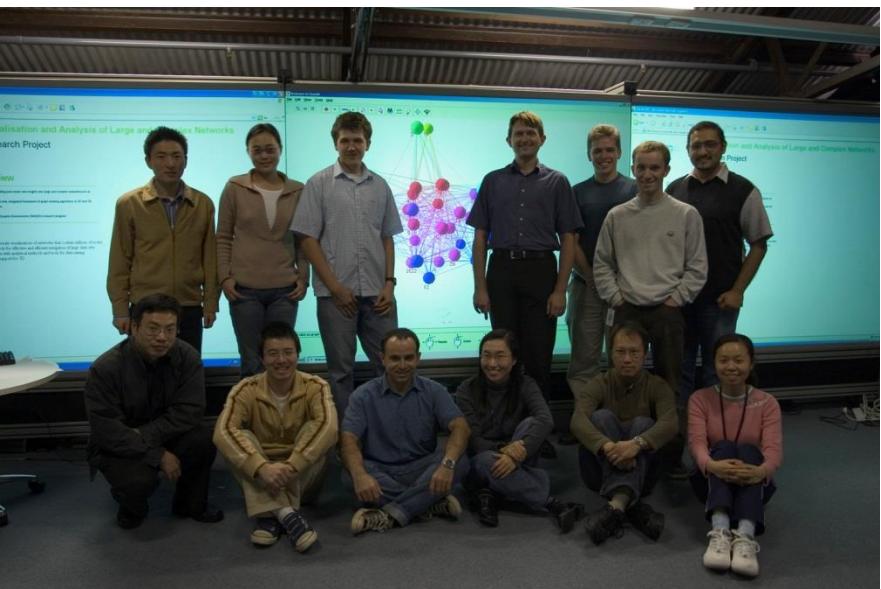
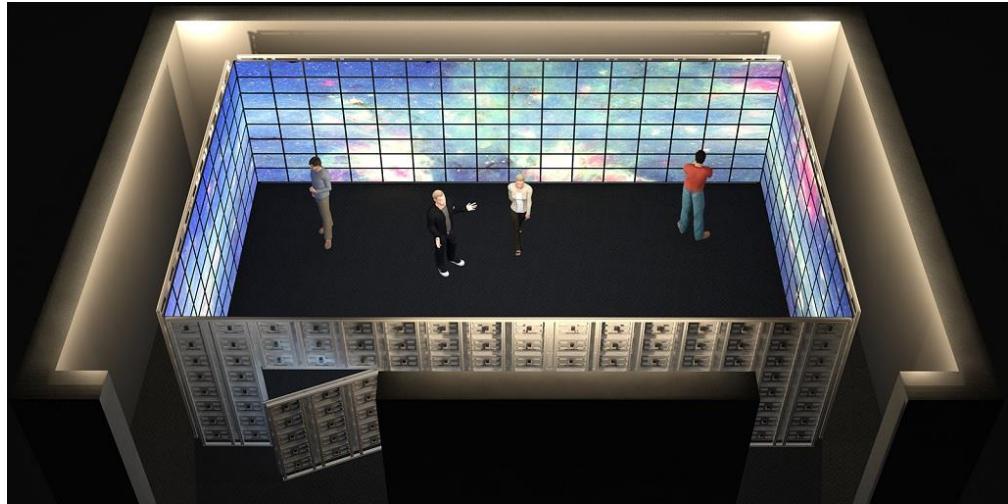
- Each key frame F_i consists of:
 - a) A queue of focus nodes
 - b) Nodes connected to a focus node.



Interaction can help two problems:

- Computational complexity: the layout is only computed for the key frame (a relatively small network)
- Visual complexity: At any one time, only the key frame (a small network) is on the screen

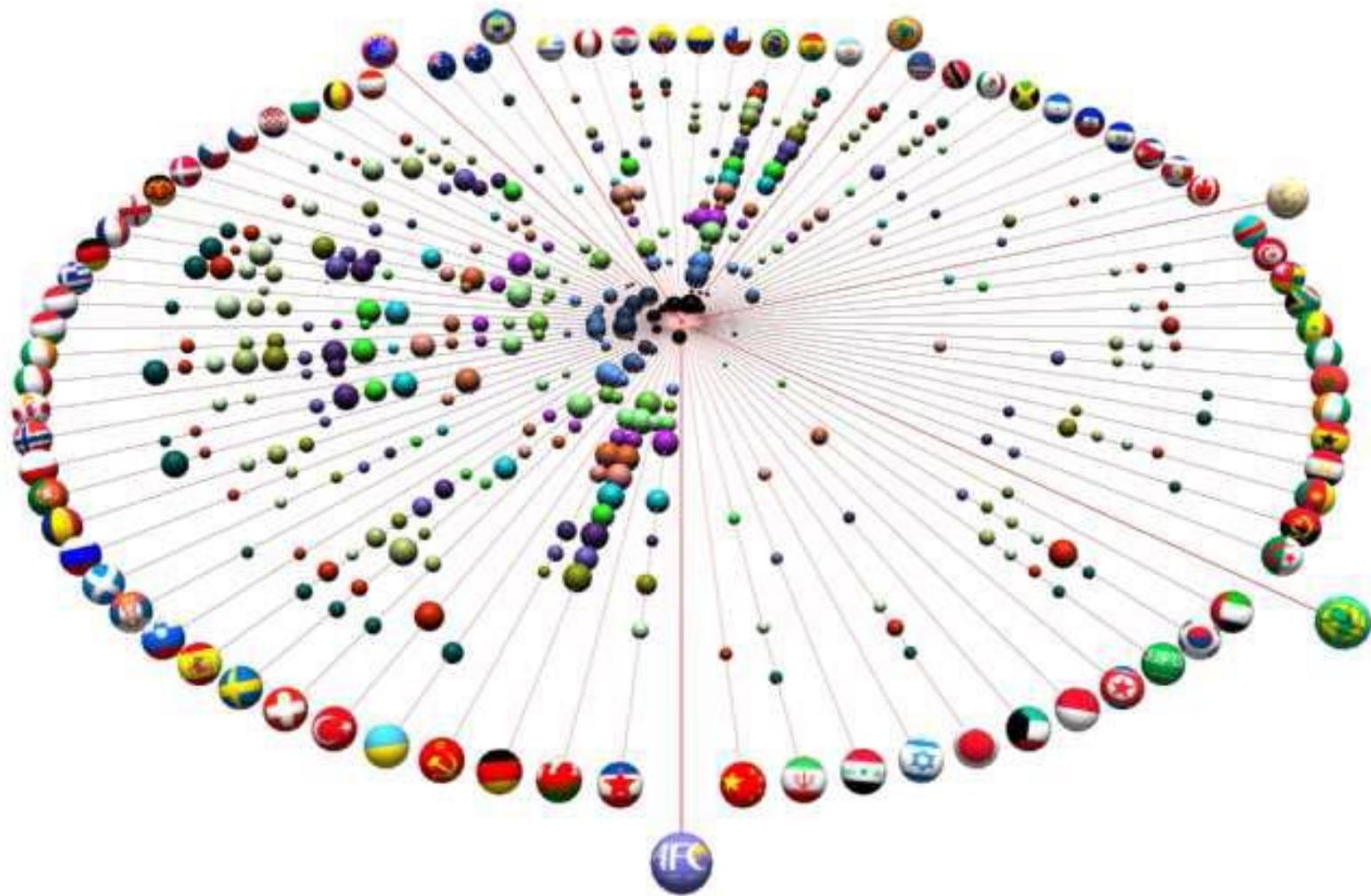
Powerwall: large screen tiled HD display



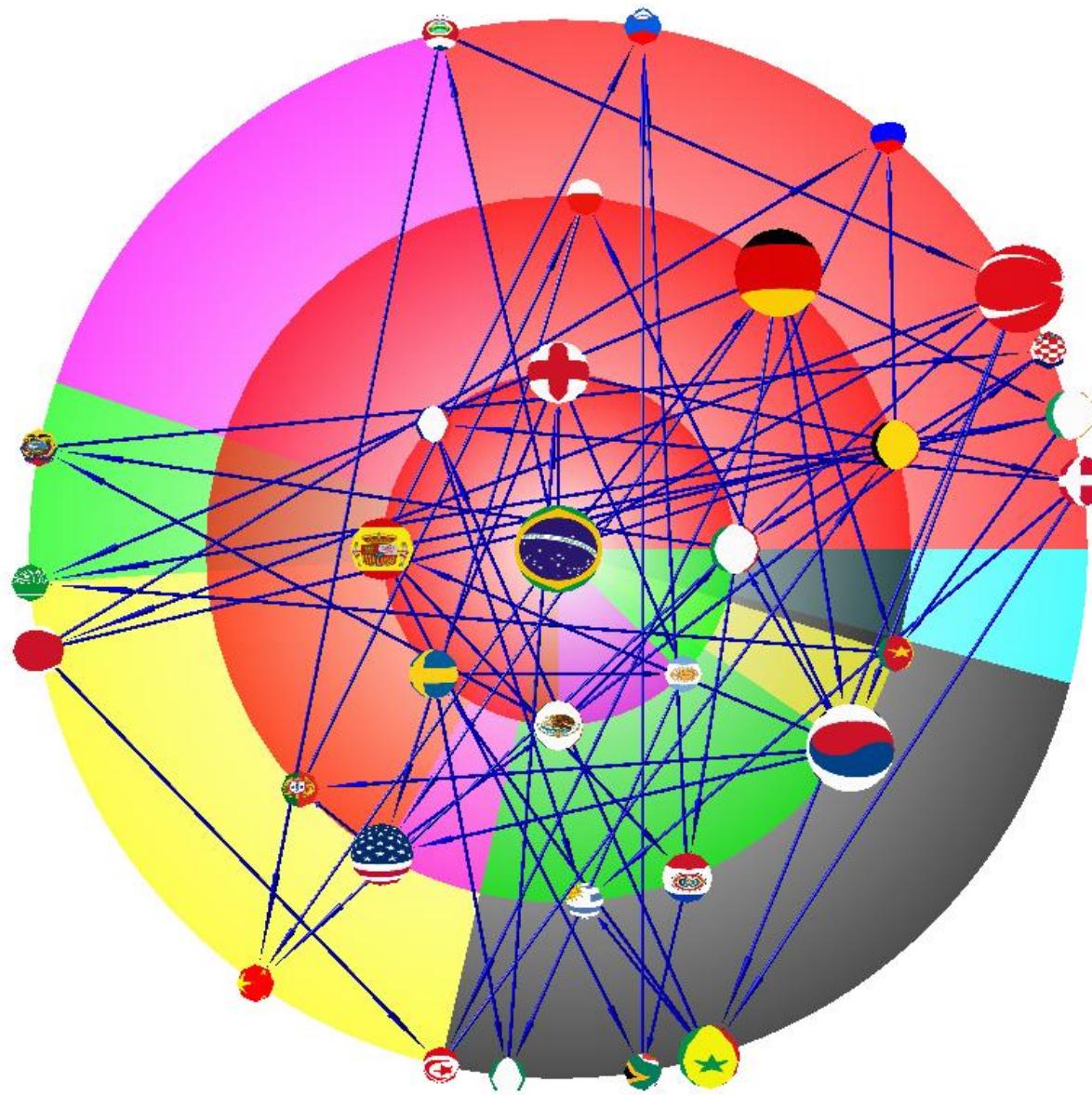
Example:

**Integration of Network Analysis
with Visualisation**

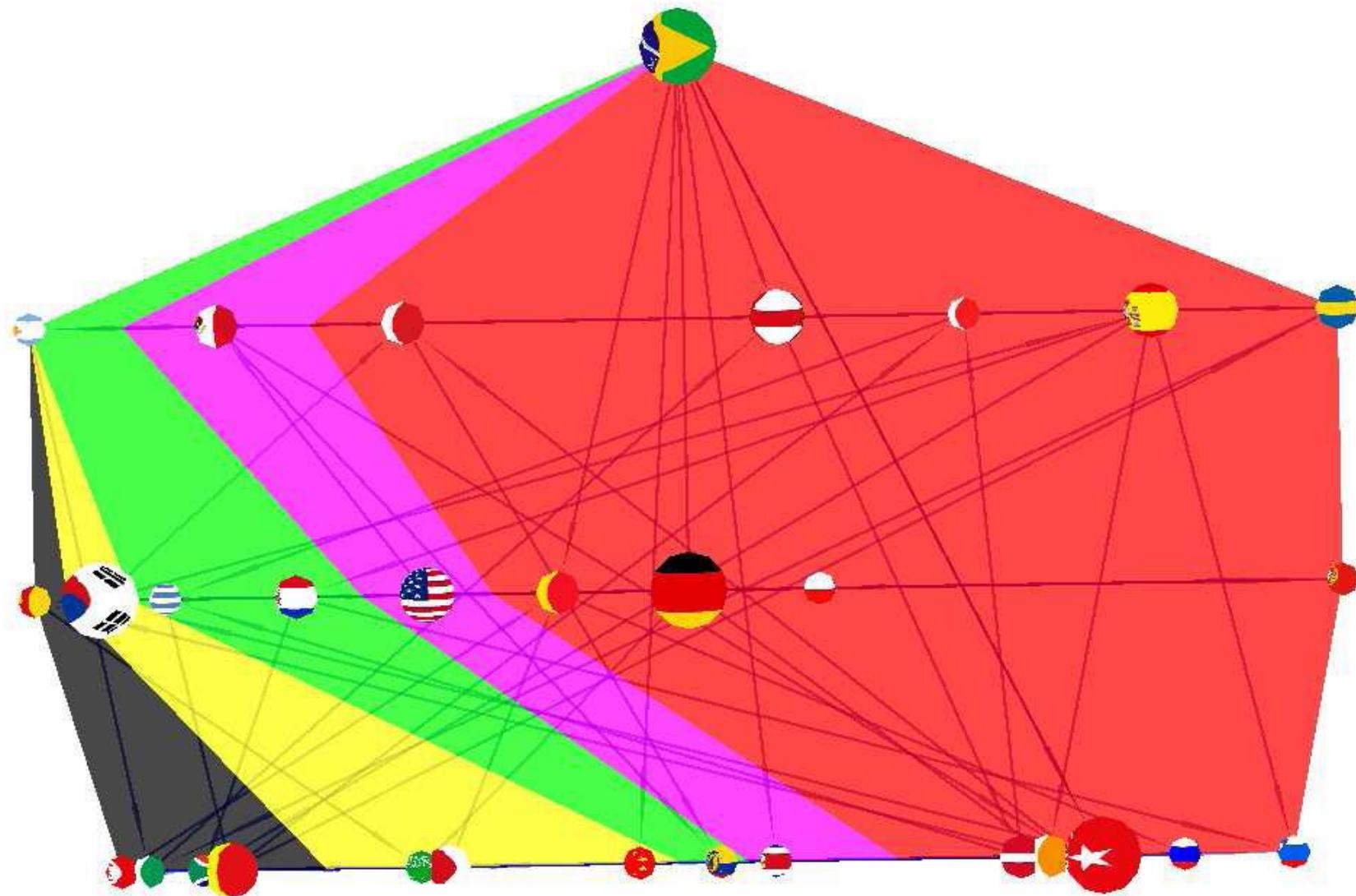
History of World Cup



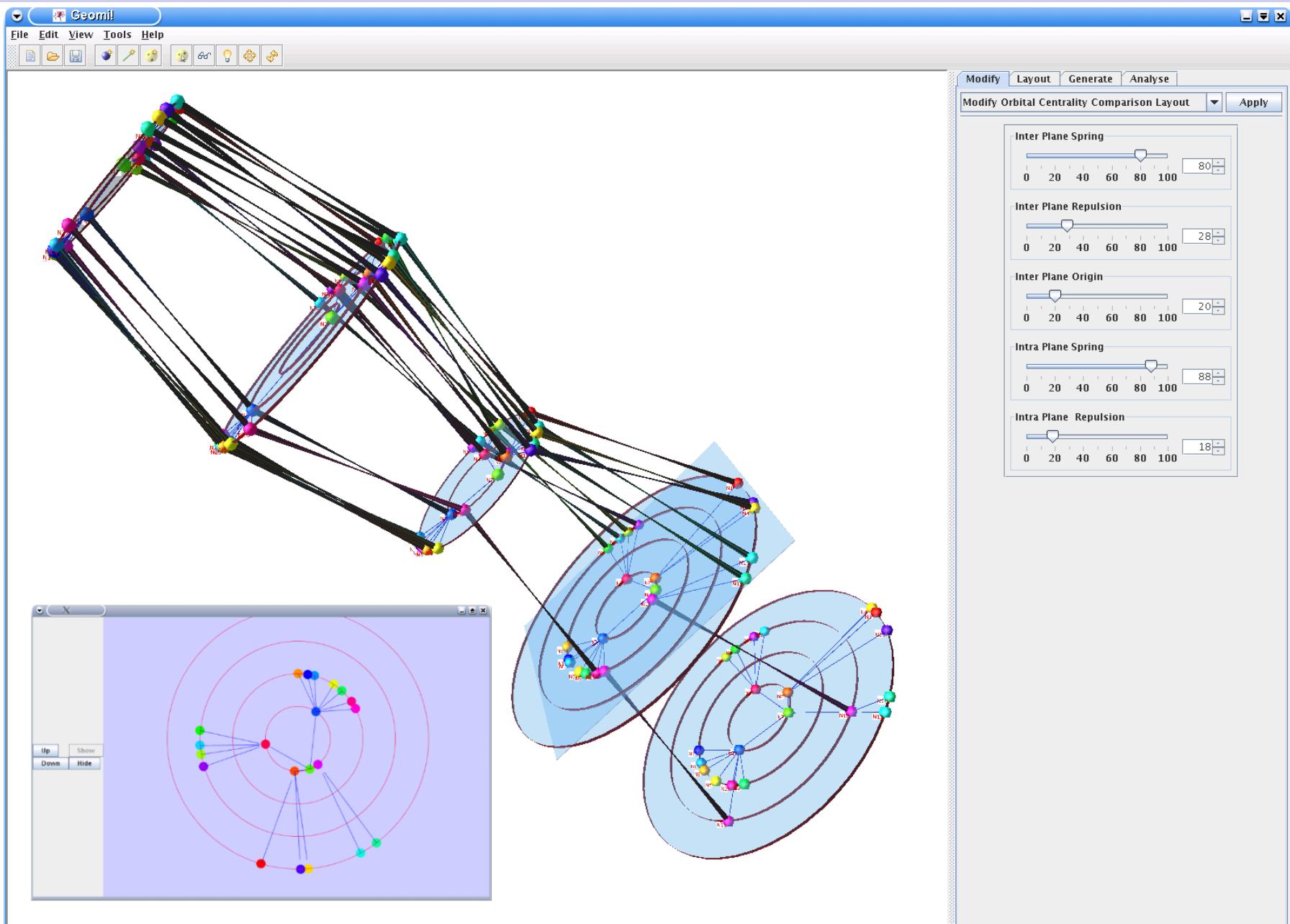
World Cup 2002



2002

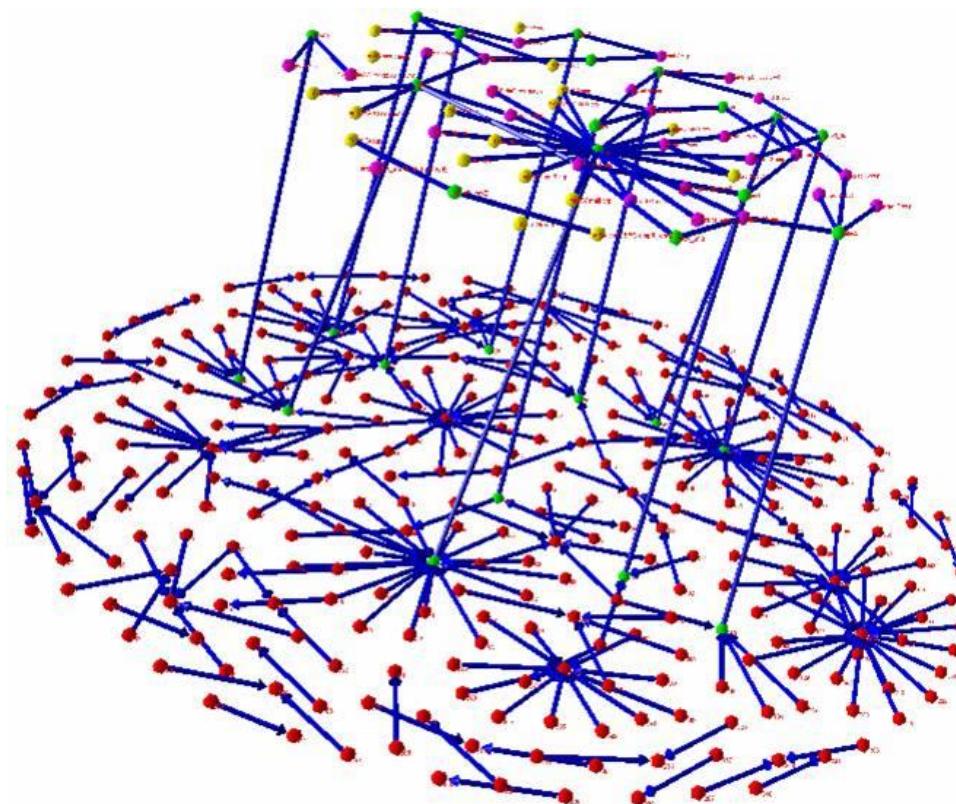


Visual Comparison of Network Centralities



Visualisation of Patterns – Network Motif

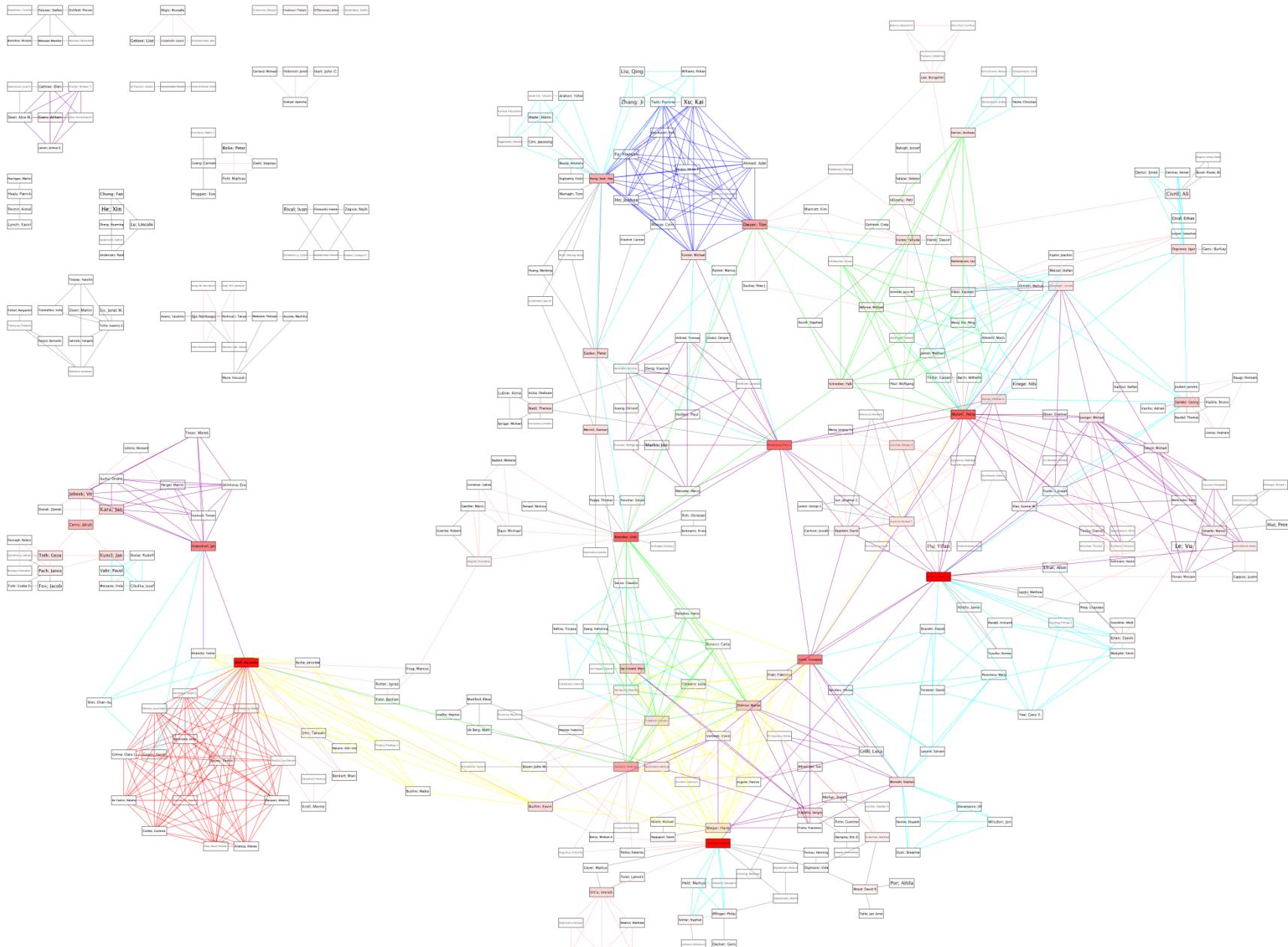
- “Motifs”: small pattern in the network which occurs with significantly high frequency.
- Data: transcriptional regulation network of Escherichia coli.

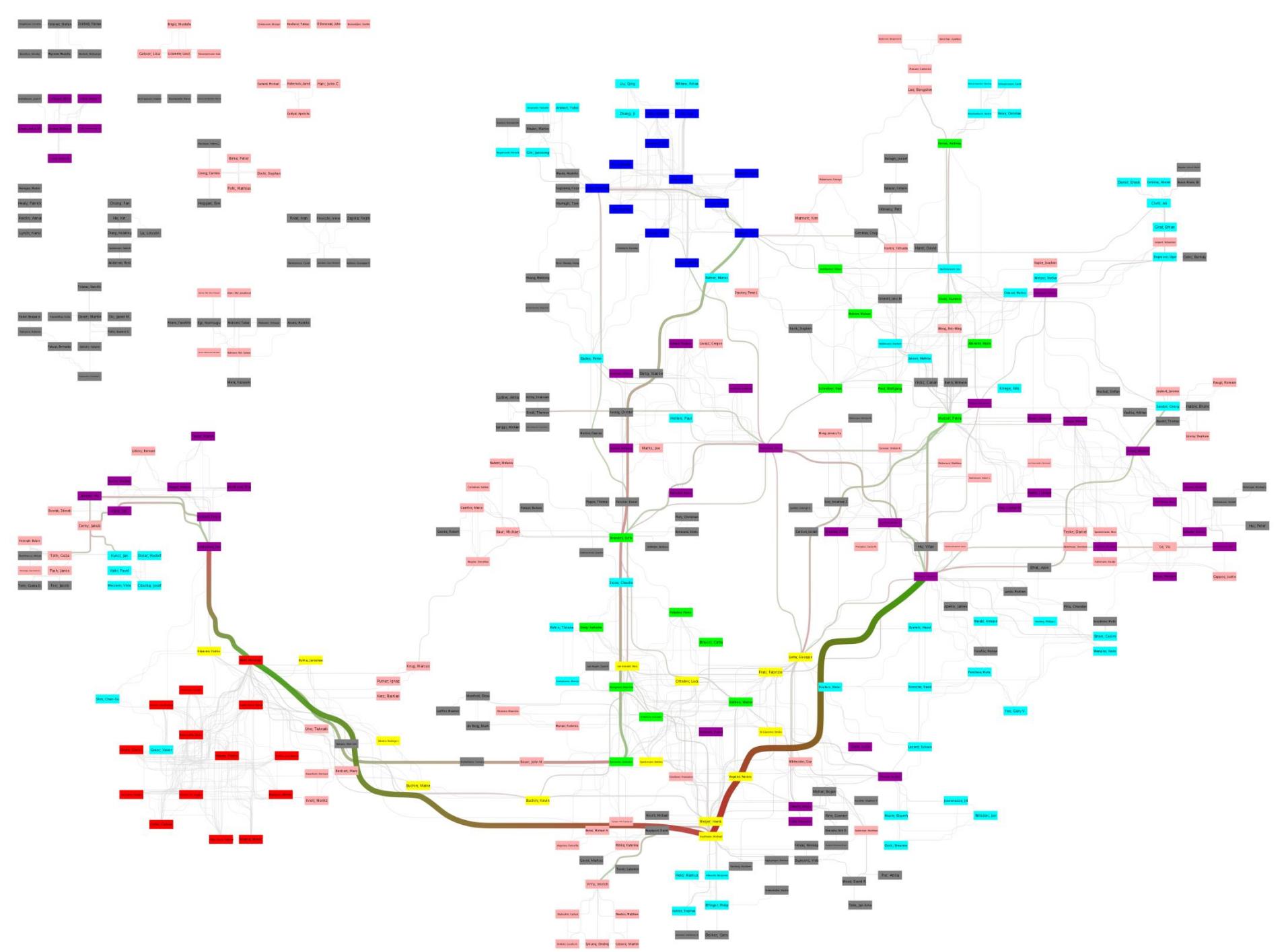


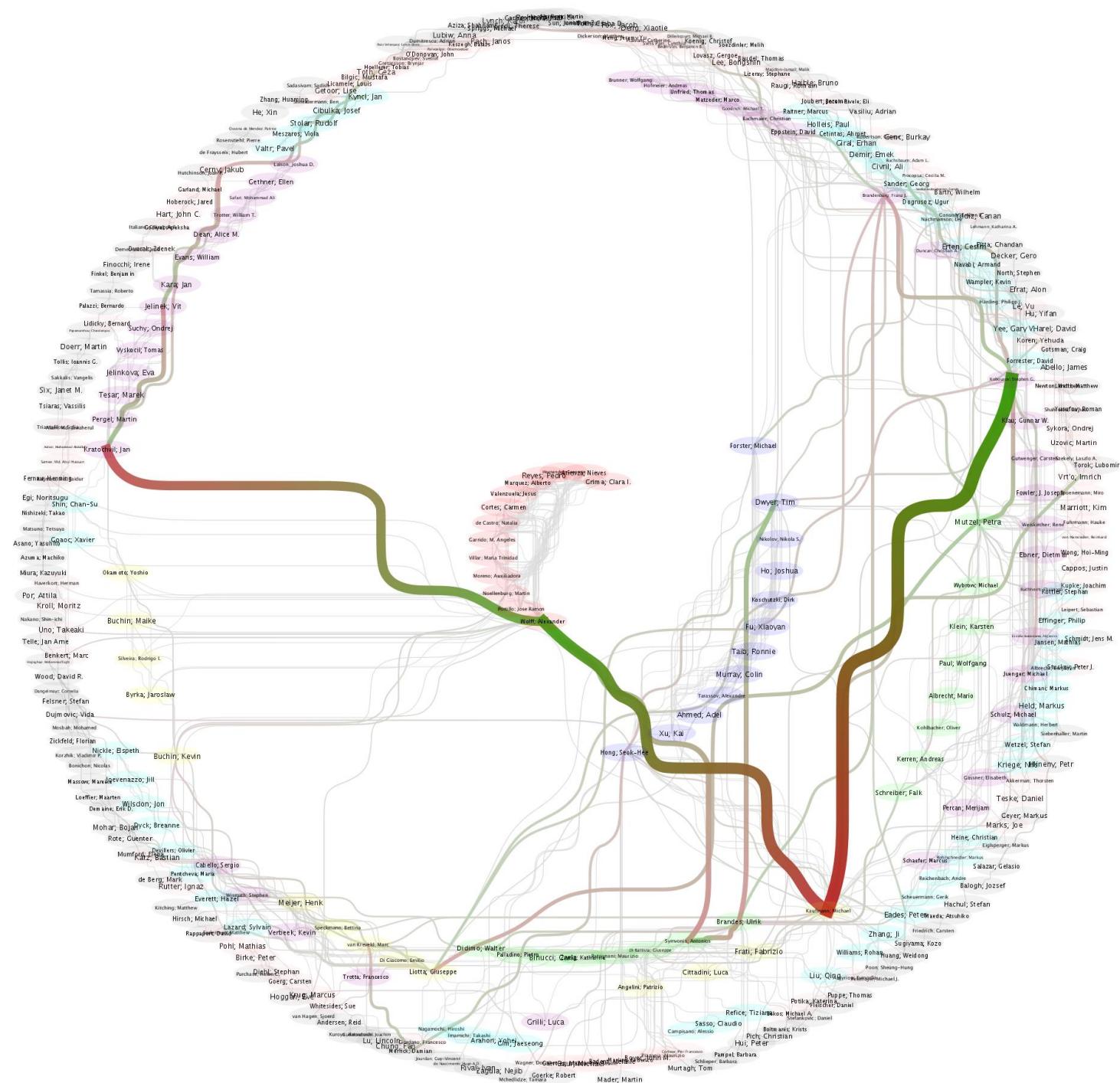
Example:

**Reducing Visual Complexity
with Edge Bundling
combined with Analysis**

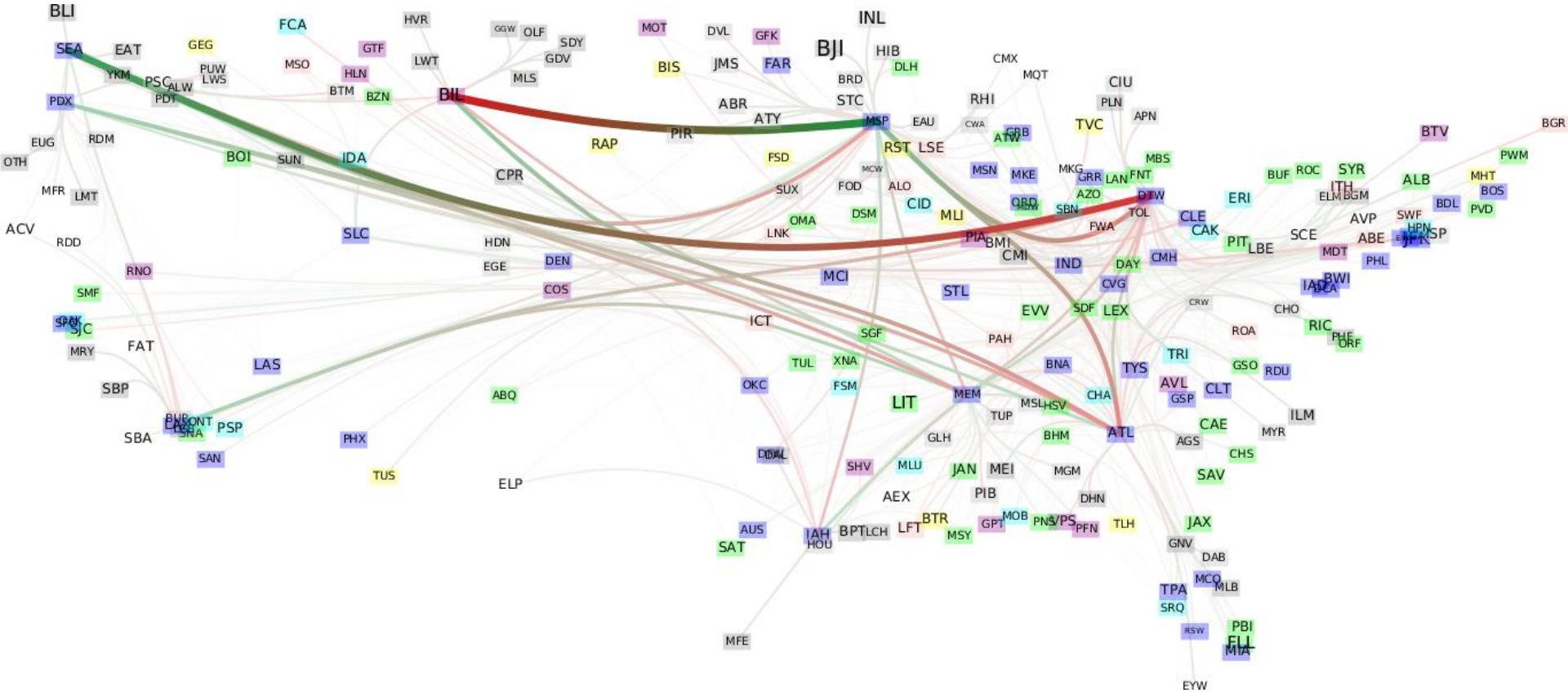
GD Research Collaboration Network

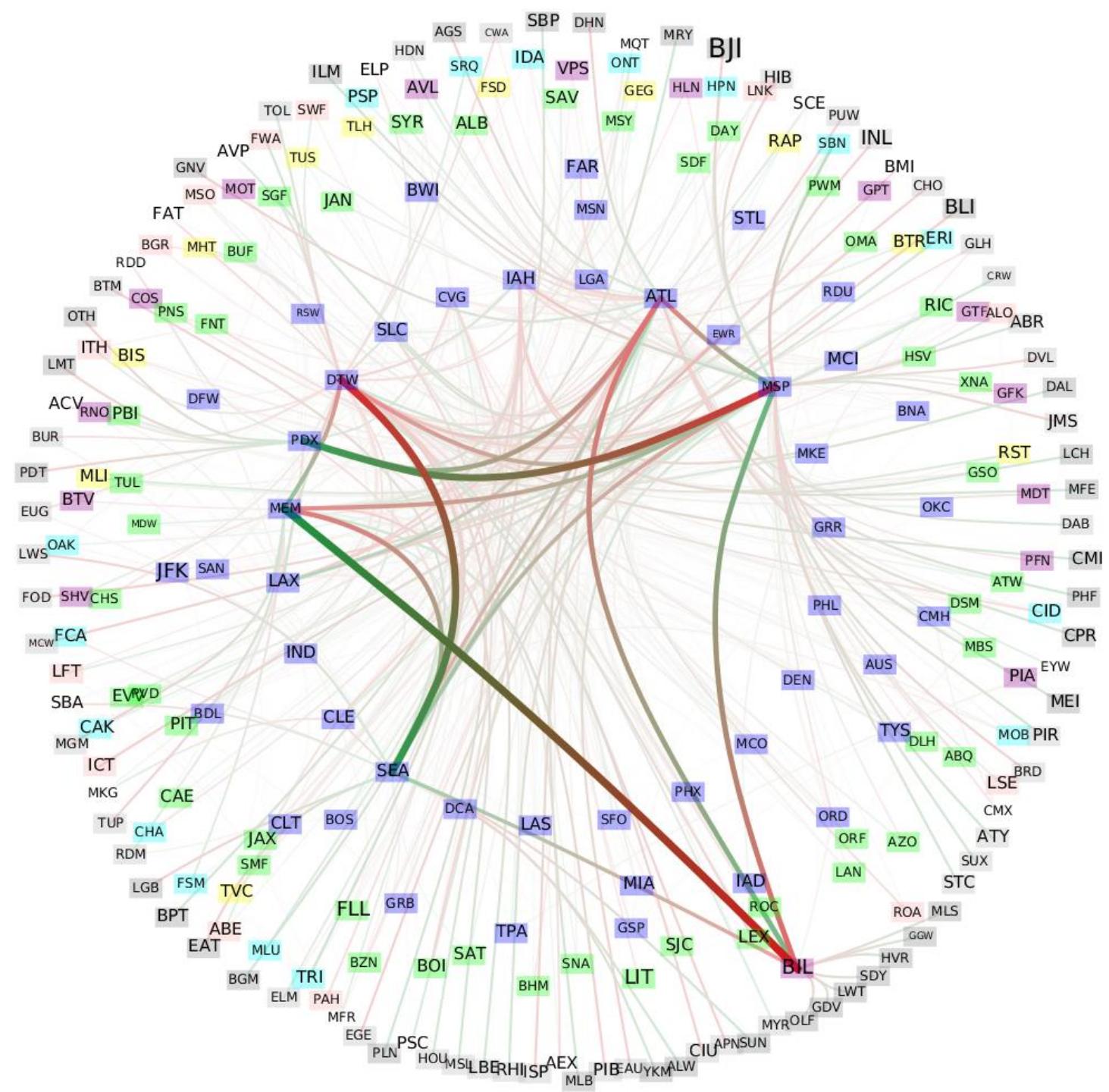




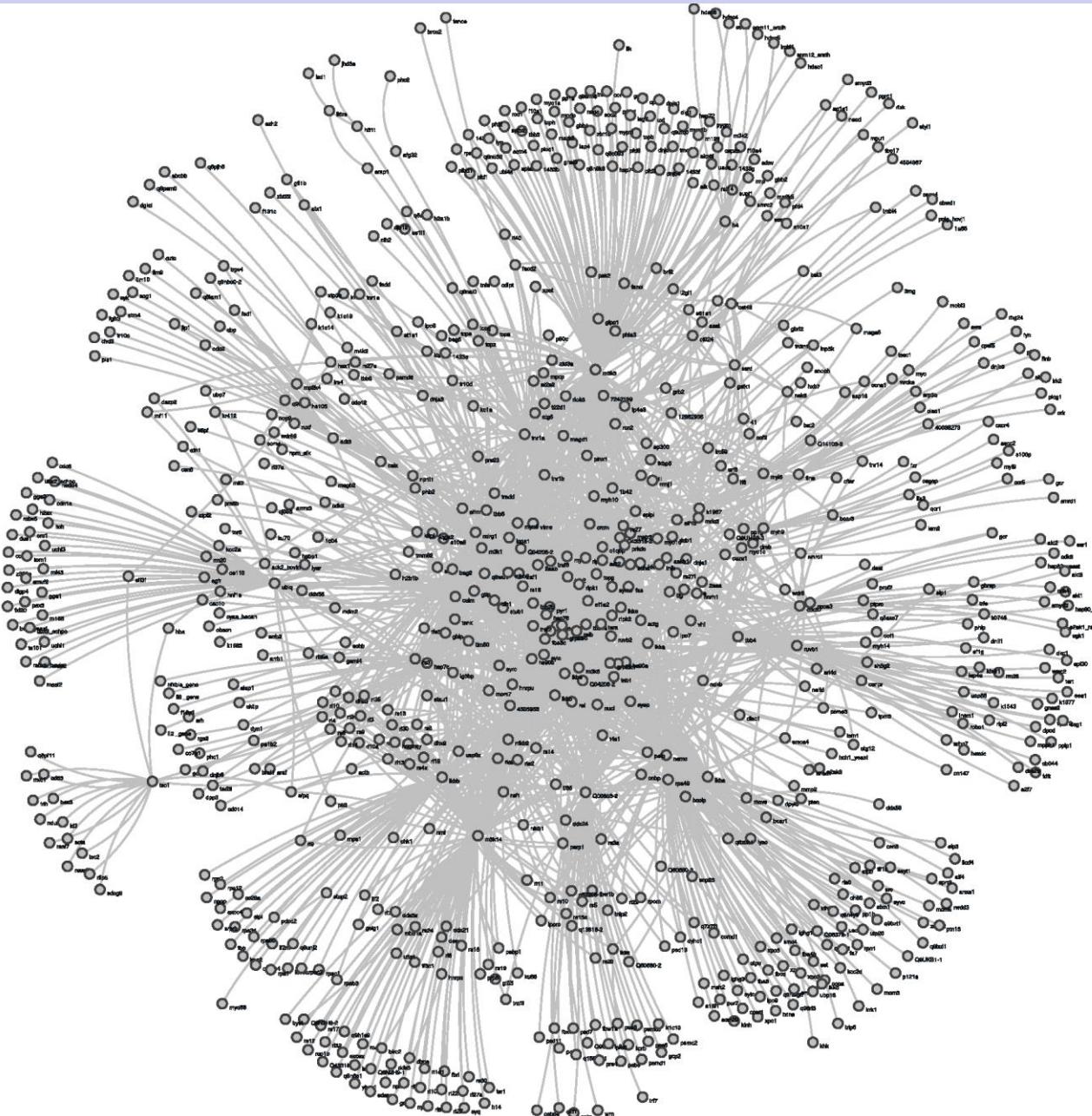


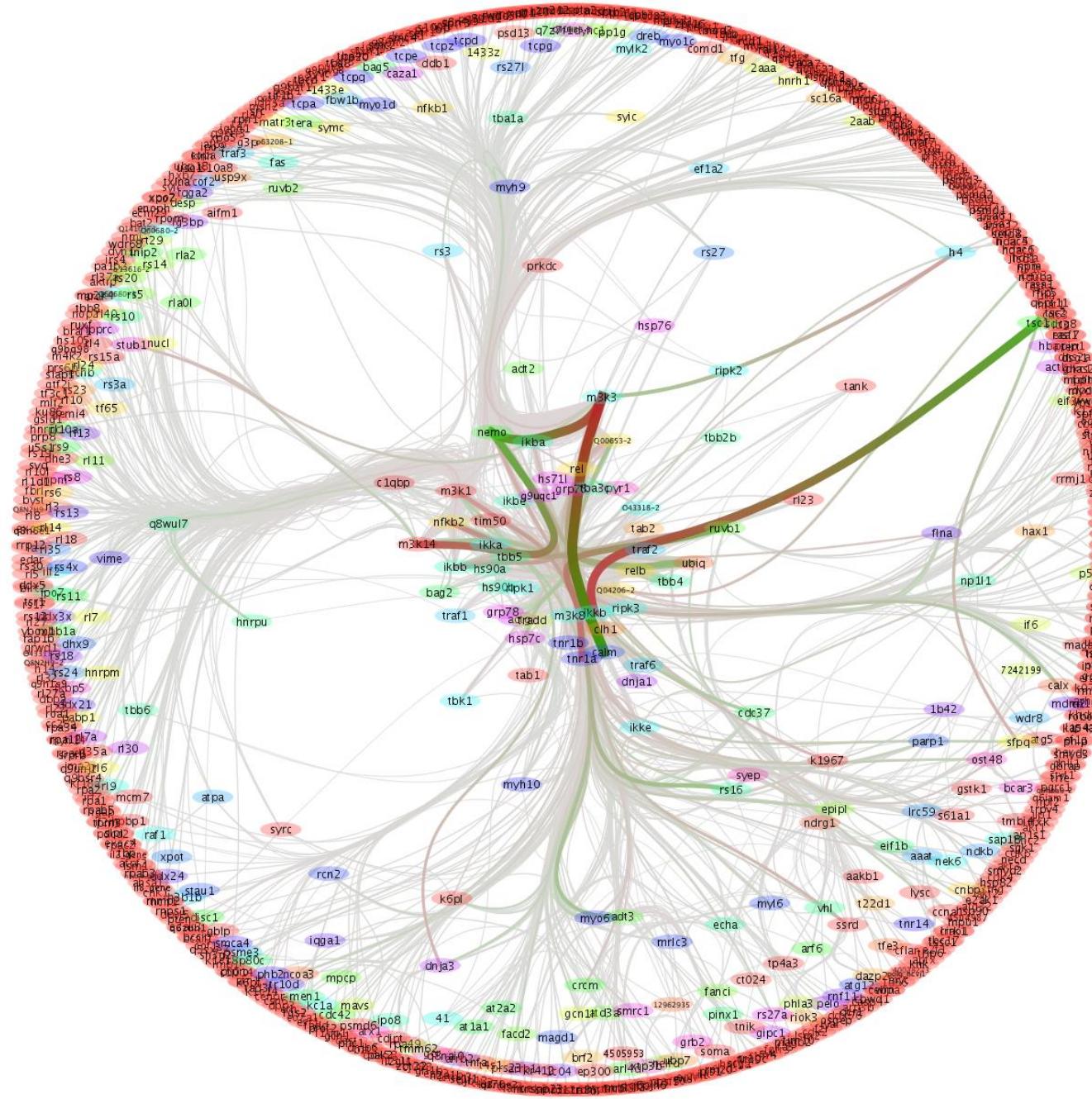
US Flight Network Traffic Analysis





PPI Network





Summary

Main Approaches for Big Data Visualisation:

1. **Cluster the data**
2. **Multi-level approach**
3. **Use 3 dimensions**
4. **Reduce Visual Complexity**
5. **Integration with Analysis**
6. **Integration with Interaction**