

VoG: Summarizing and Understanding Large Graphs

Danai Koutra

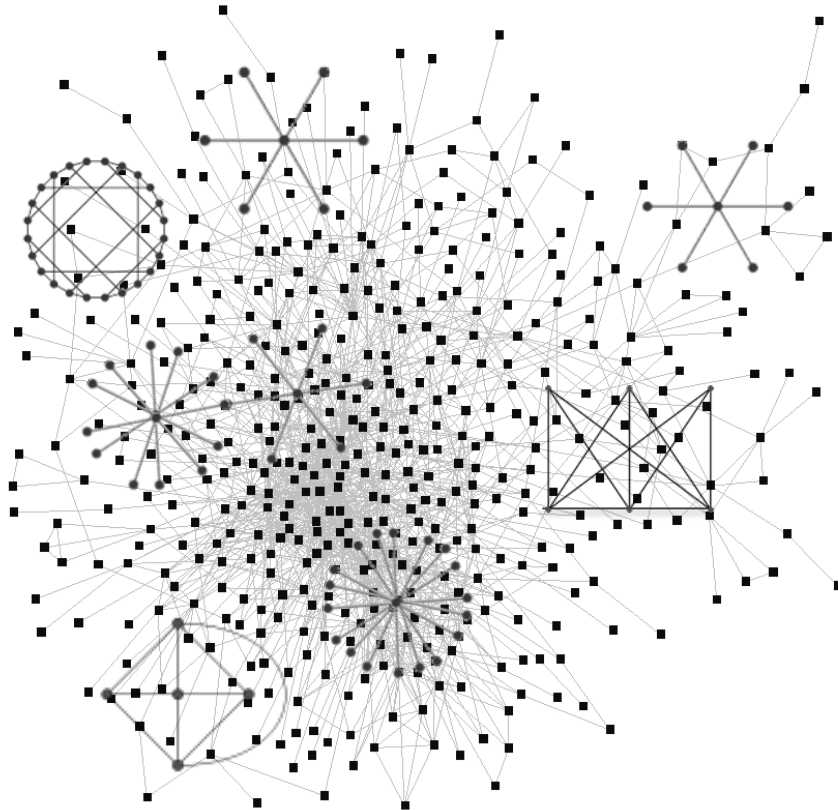
Jilles Vreeken

U Kang

Christos Faloutsos

SDM, 24-26 April 2014, Philadelphia, USA

Problem Definition: Graph Summarization



Given: a graph

Find: a succinct summary
with possibly
overlapping subgraphs

\approx

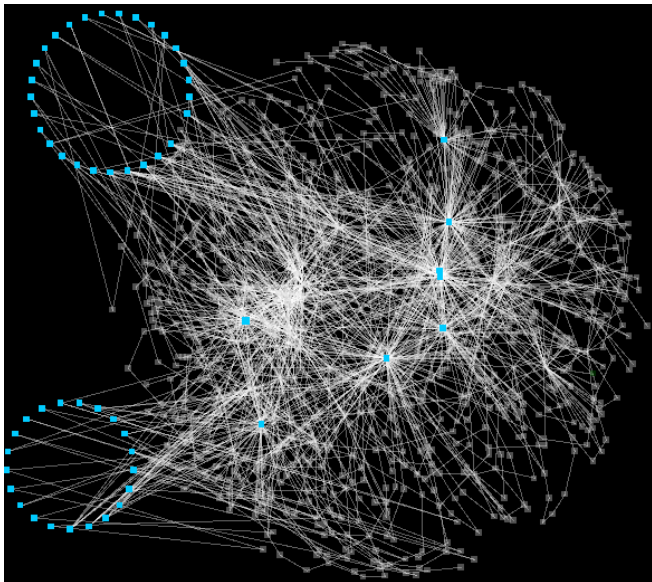
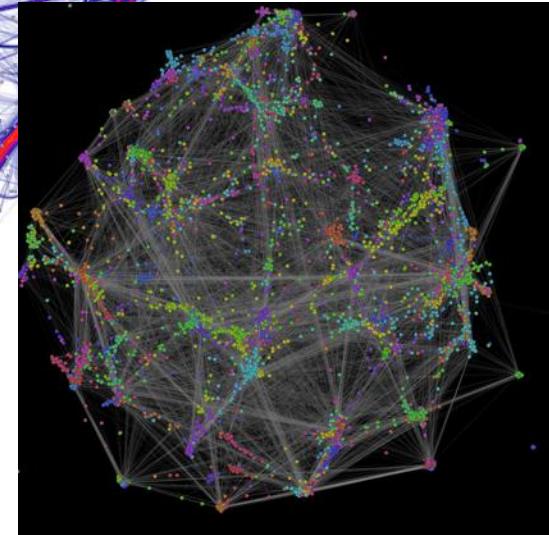
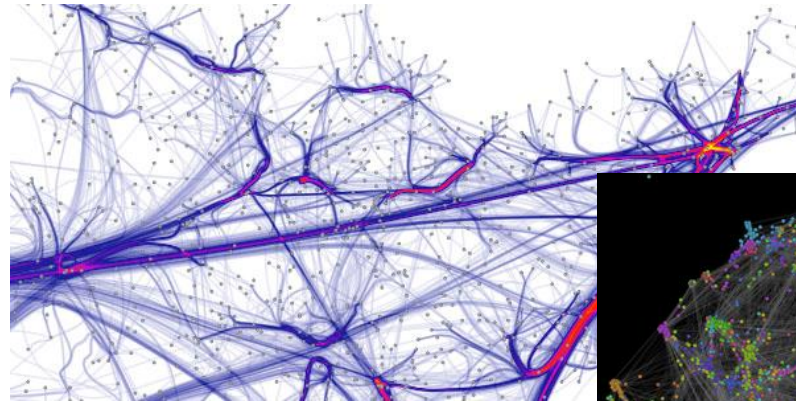
*important graph
structures*



Why graph summarization?



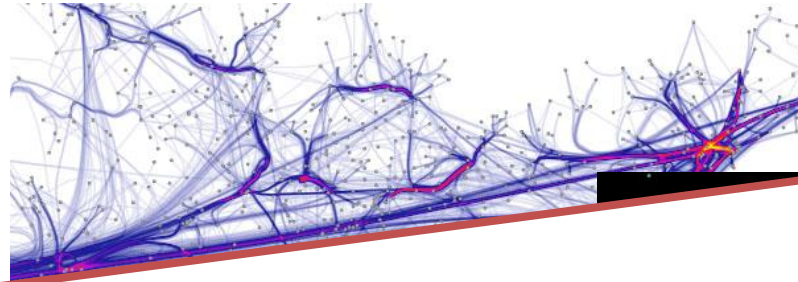
Visualization



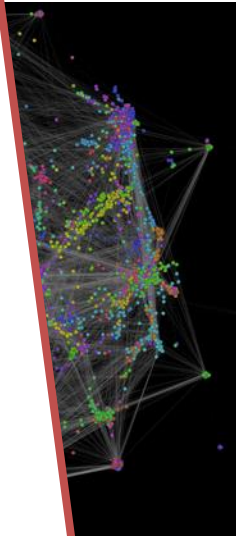
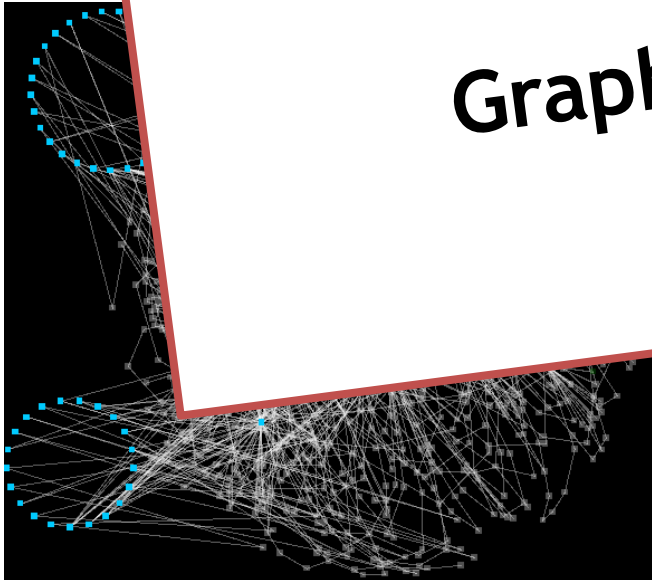
Guiding attention



Why graph summarization?



Graph Understanding



Roadmap

Main Idea

Proposed Algorithm: VoG

VoG: Step-by-Step

Experiments

Conclusions



Main Idea



1) Use a graph vocabulary:



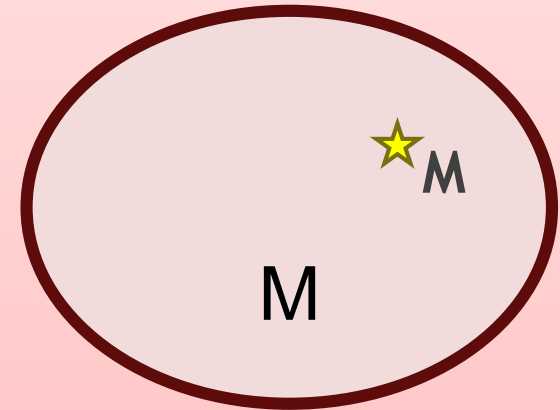
2) Best graph summary

→ structures with the lowest
compression cost (MDL)



Minimum Description Length Principle

- Given a set of models M ,
- the best model $M \in M$ is



$$\operatorname{argmin}_M L(M) + L(D|M)$$

bits
for M

bits for the
data using M

Formally: Minimum Graph Description

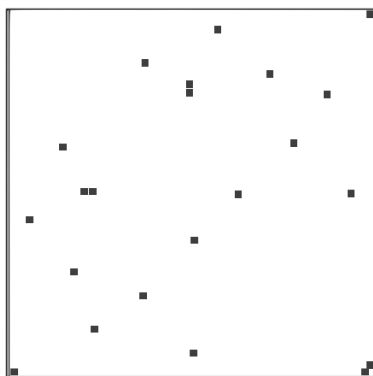


Given: - a graph G with adjacency matrix A
- vocabulary Ω

Find: model M s.t.

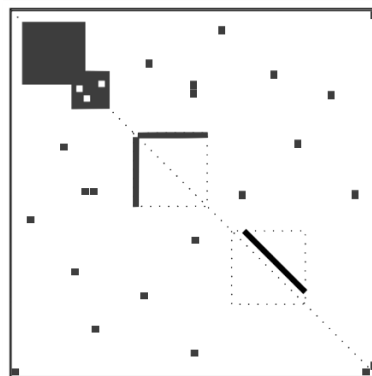
$$\min L(G, M) = \min L(M) + L(E)$$

Error E



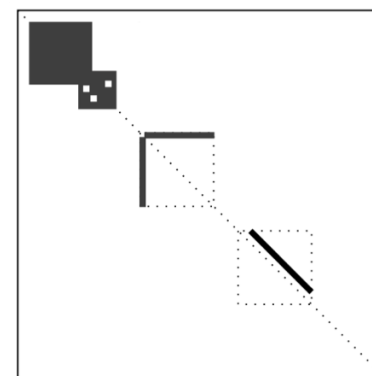
=

Adjacency A



\otimes

Model M



Roadmap

Main Idea

Proposed Algorithm: VoG

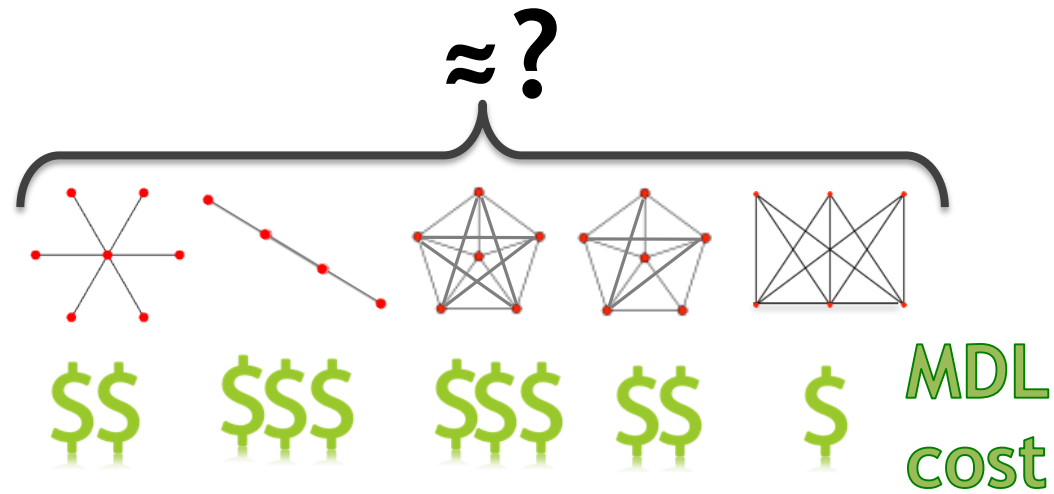
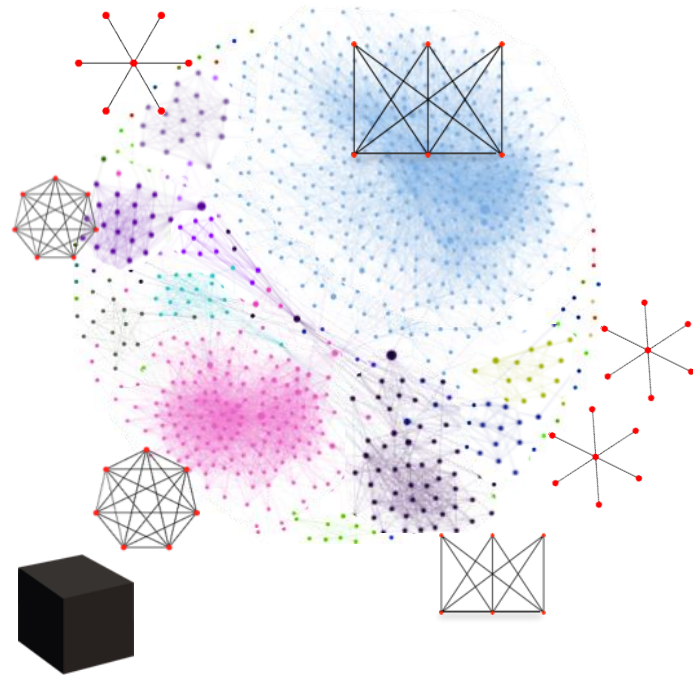
VoG: Step-by-Step

Experiments

Conclusions



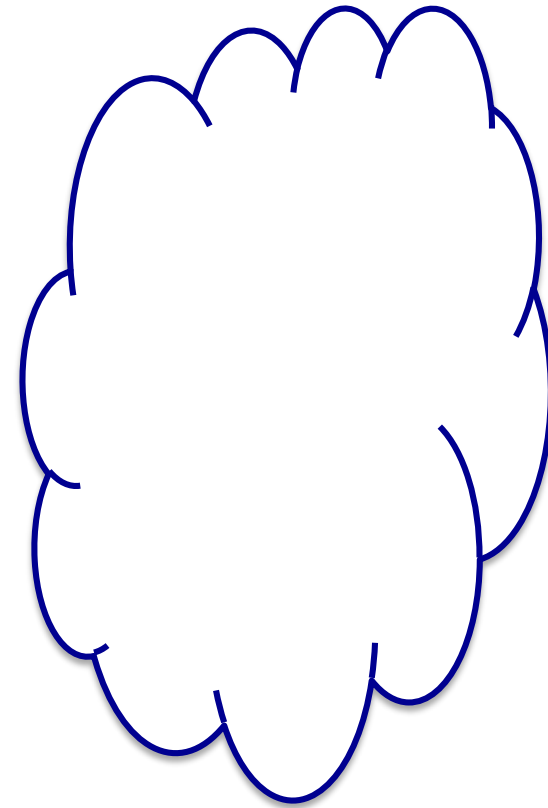
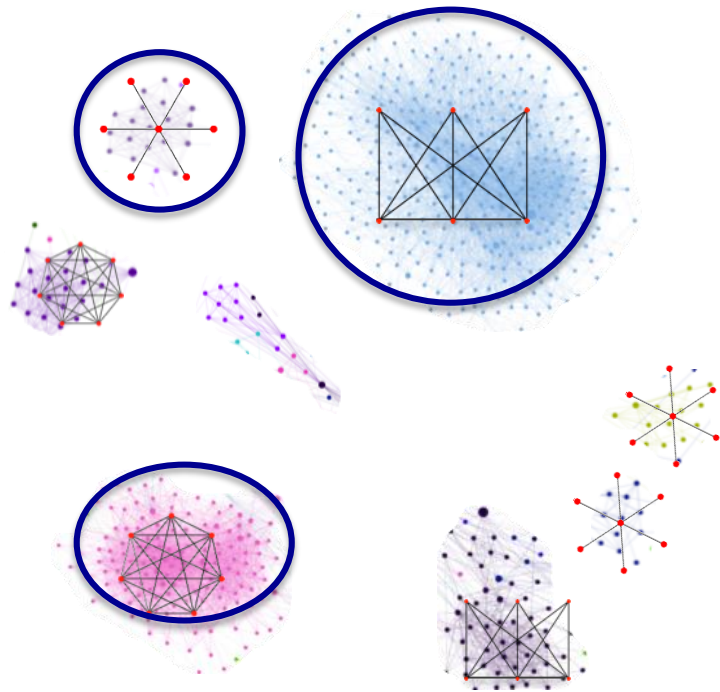
VoG: Overview



argmin



VoG: Overview



some criterion

Summary



Roadmap

Main Idea

Proposed Algorithm: VoG

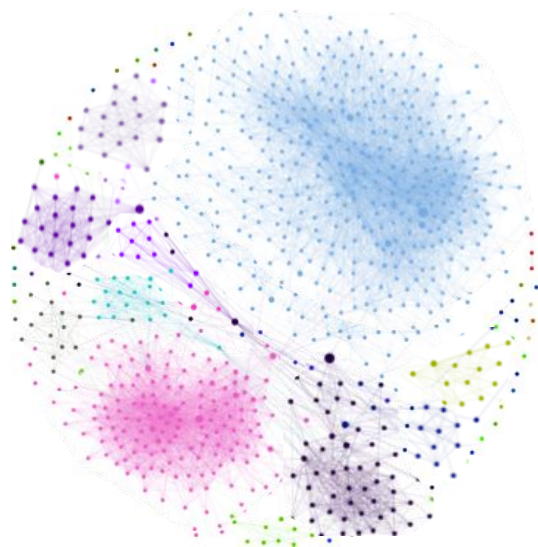
VoG: Step-by-Step

Experiments

Conclusions



Step 1: Graph Decomposition



Could use:

ANY graph decomposition
method

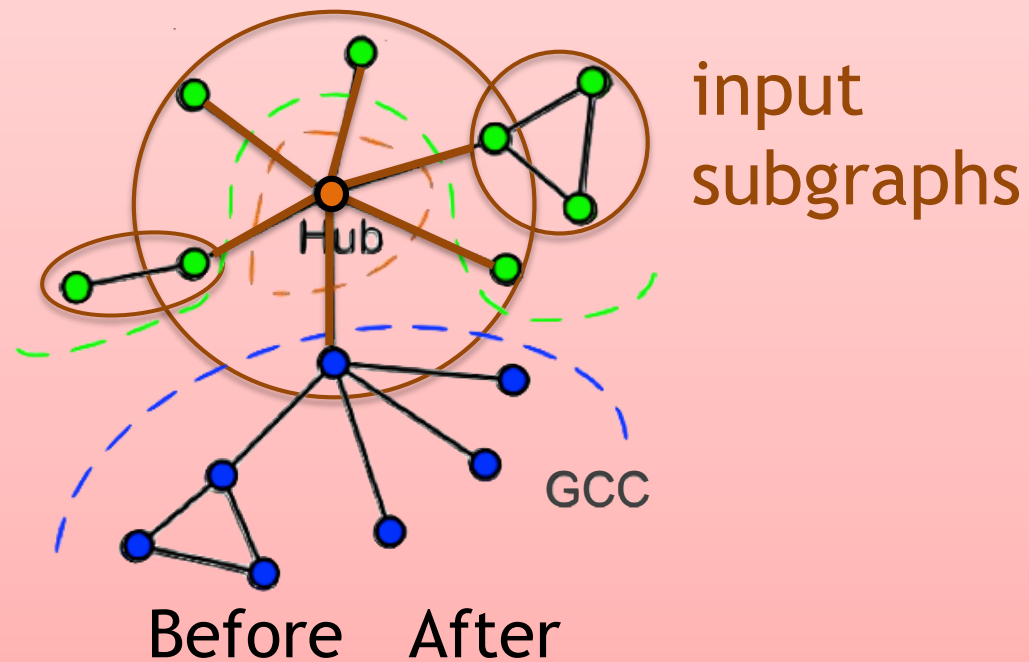
We use:

SlashBurn



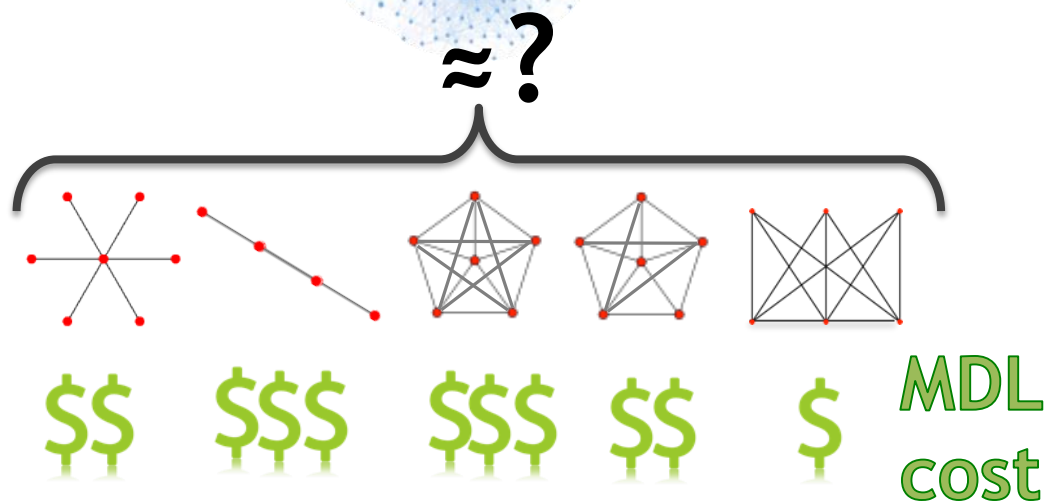
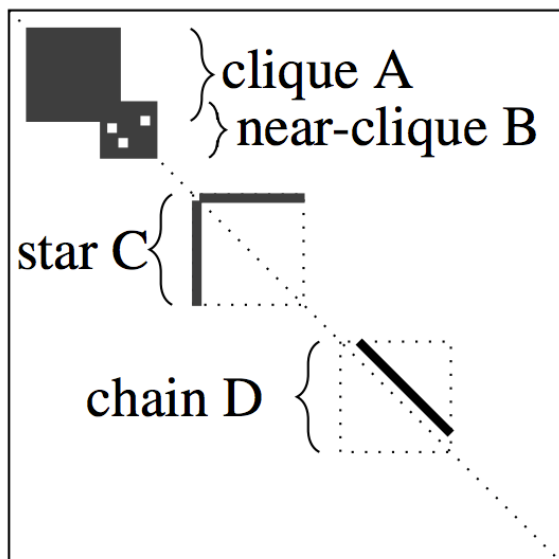
SlashBurn

- *Slash* top-k hubs, *burn* edges
- Repeat on the remaining GCC



[U Kang and Christos Faloutsos. ICDM'11]

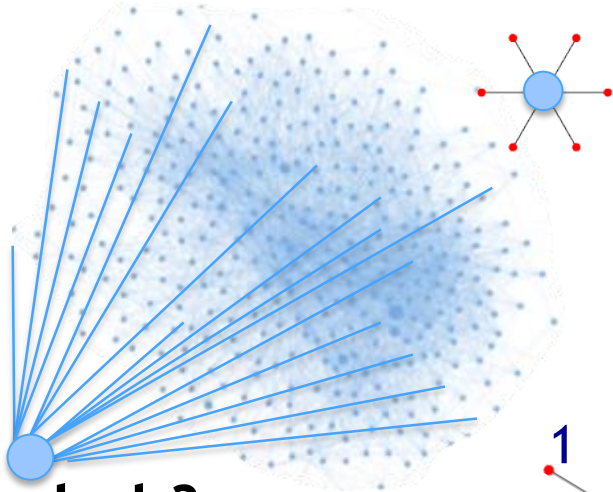
Step 2: Graph Labeling



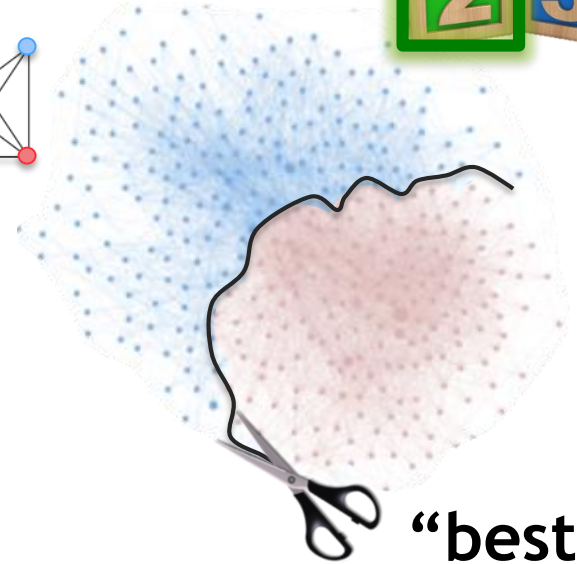
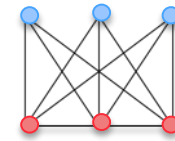
argmin



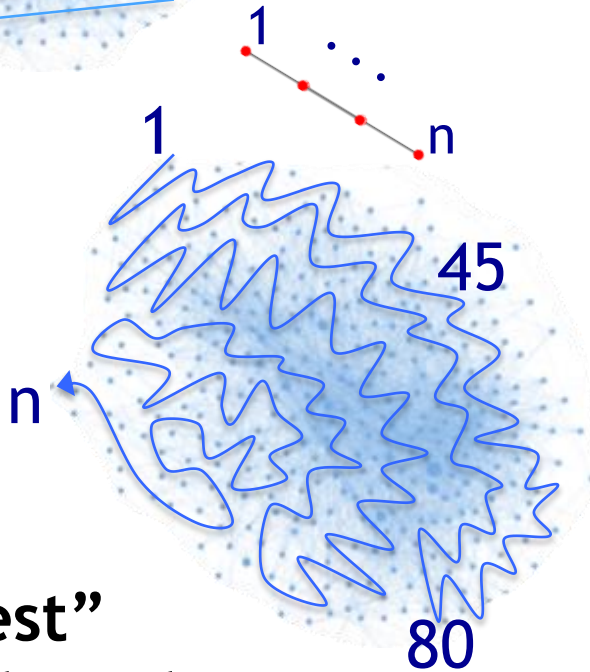
Graph Representation



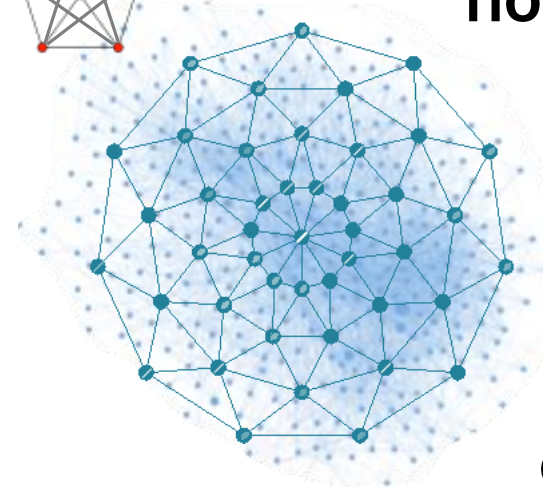
hub?



“best”
node split?



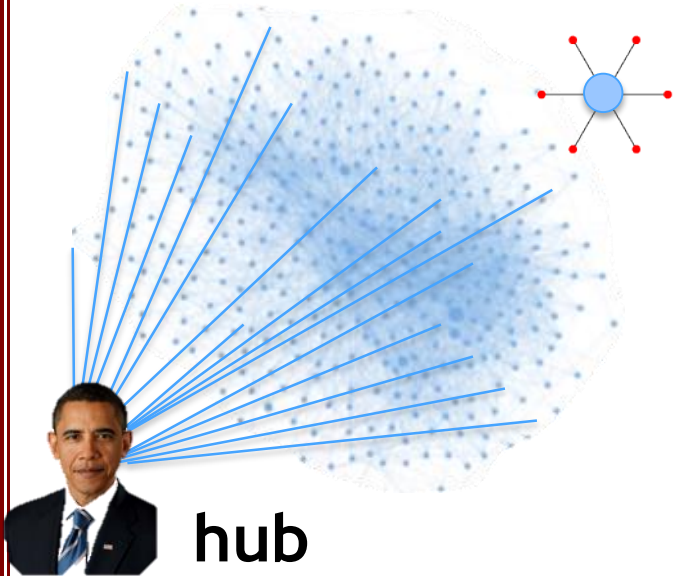
“best”
node ordering?



missing
edges?

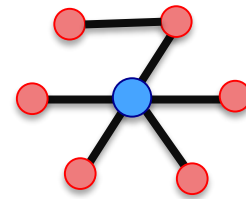


Graph Representation



Hub: top-deg node
Spokes: the rest

$n=7$



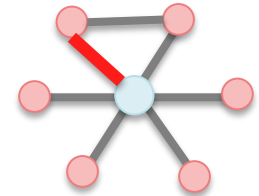
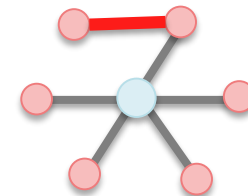
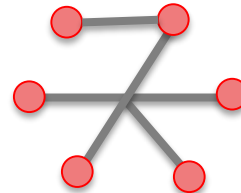
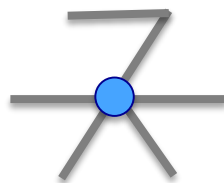
\$\$\$

Star structure

+

Errors

6



Danai Koutra - SDM'14

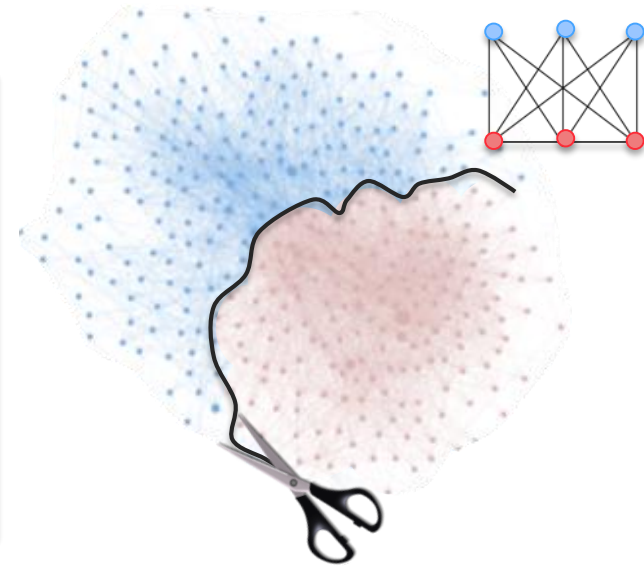
17



Graph Representation

Max bipartite graph: **NP-hard**

Heuristic: Belief Propagation with heterophily for node classification
(blue/red)

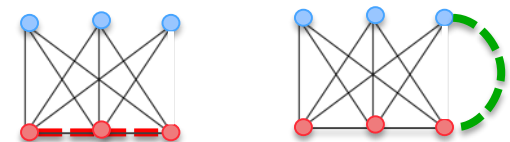


\$\$\$

Bipartite graph structure

+

Errors

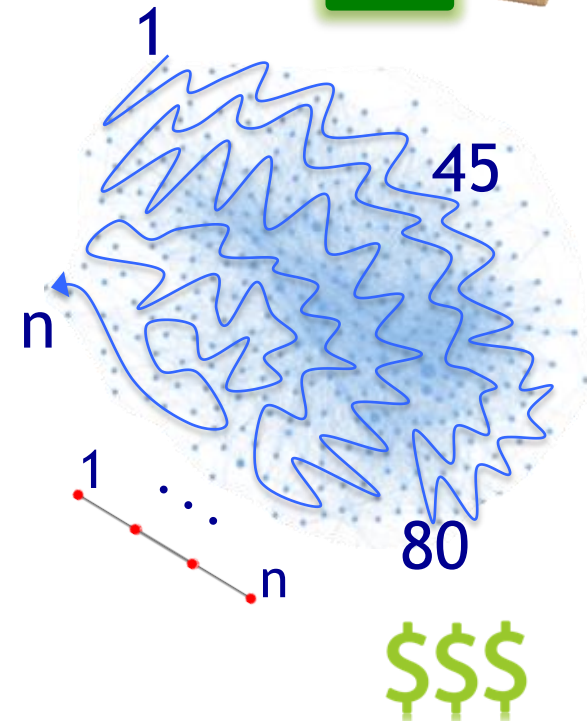


Graph Representation



Longest path: **NP-hard**

Heuristic: BFS + local search



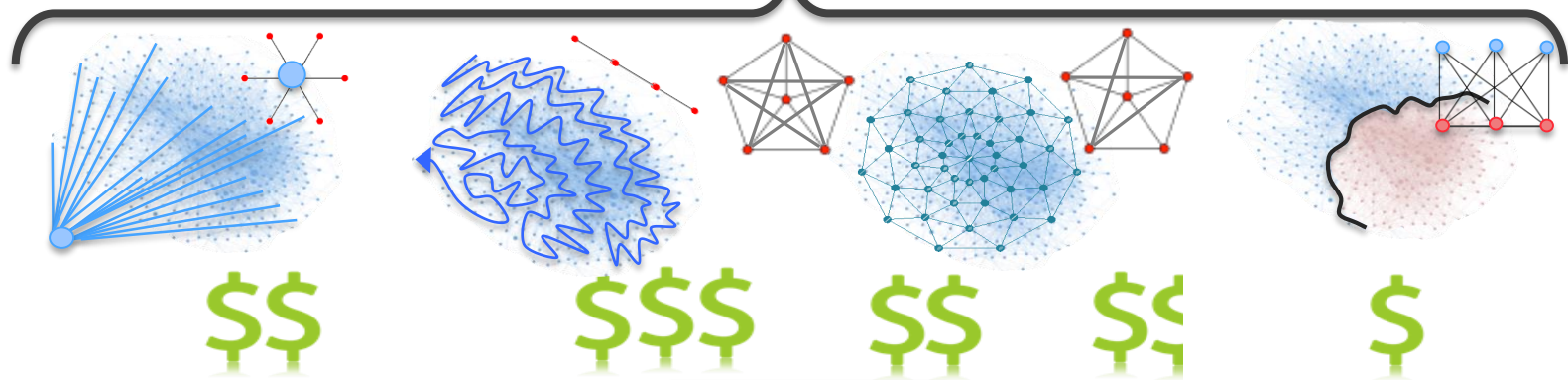
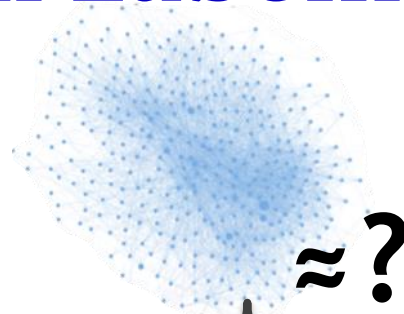
Chain structure

+

Errors



Step 2: Graph Labeling



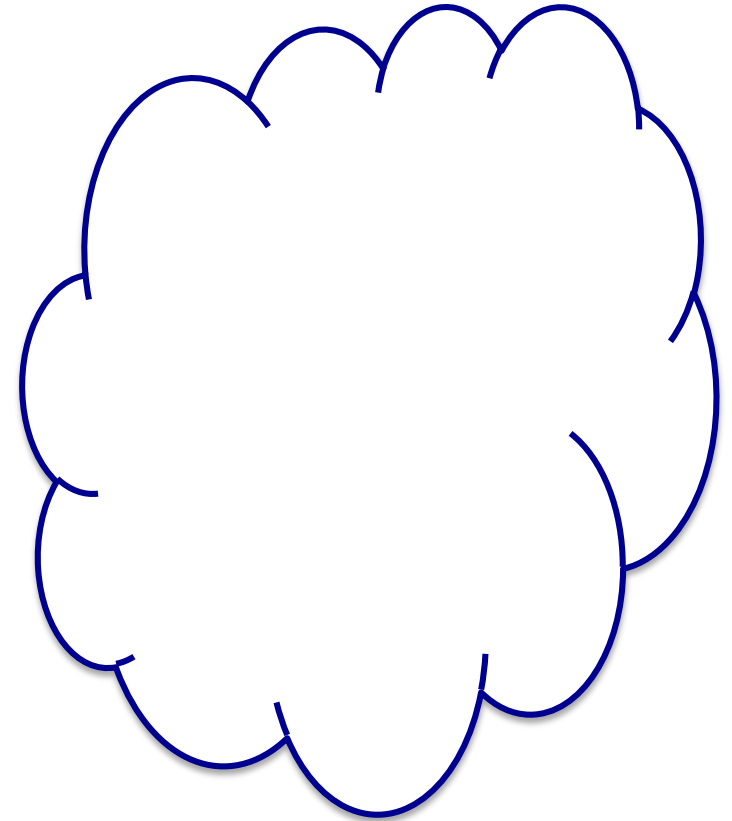
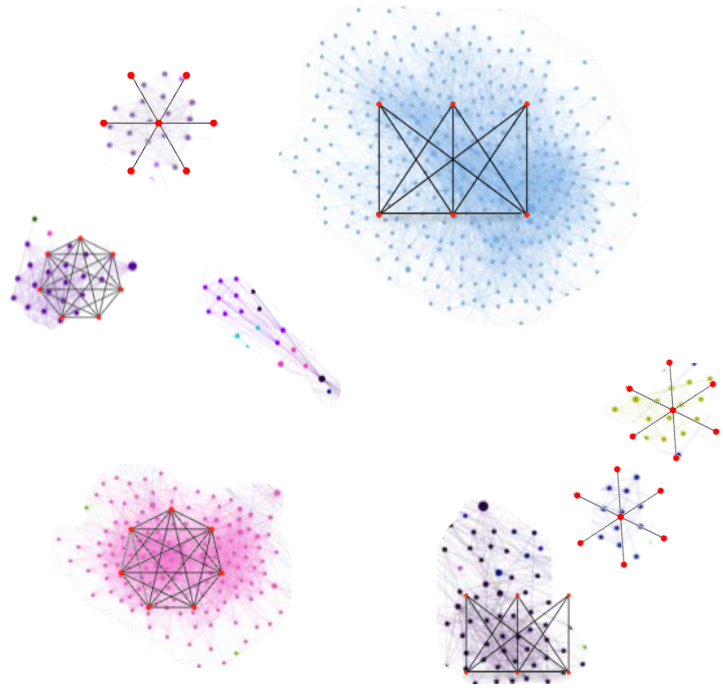
MDL
Cost
 $L(m) + L(e)$



Step 3: Summary Assembly

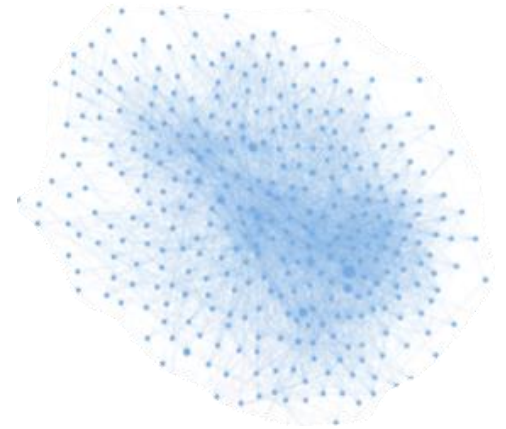
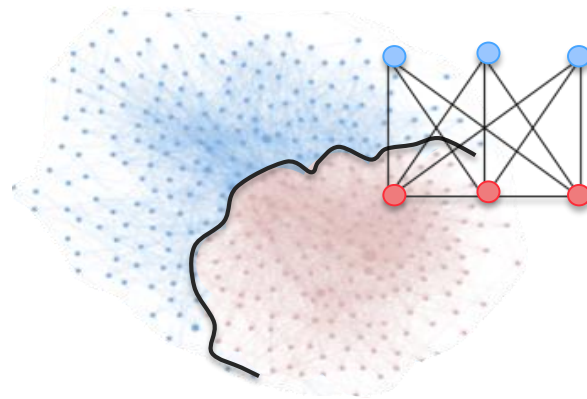
PLAIN

Summary



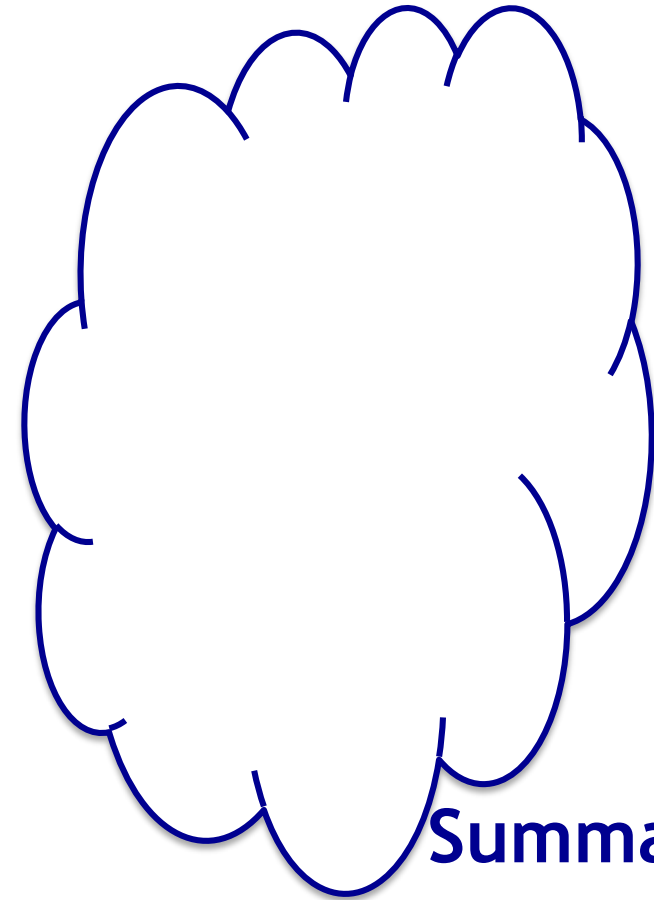
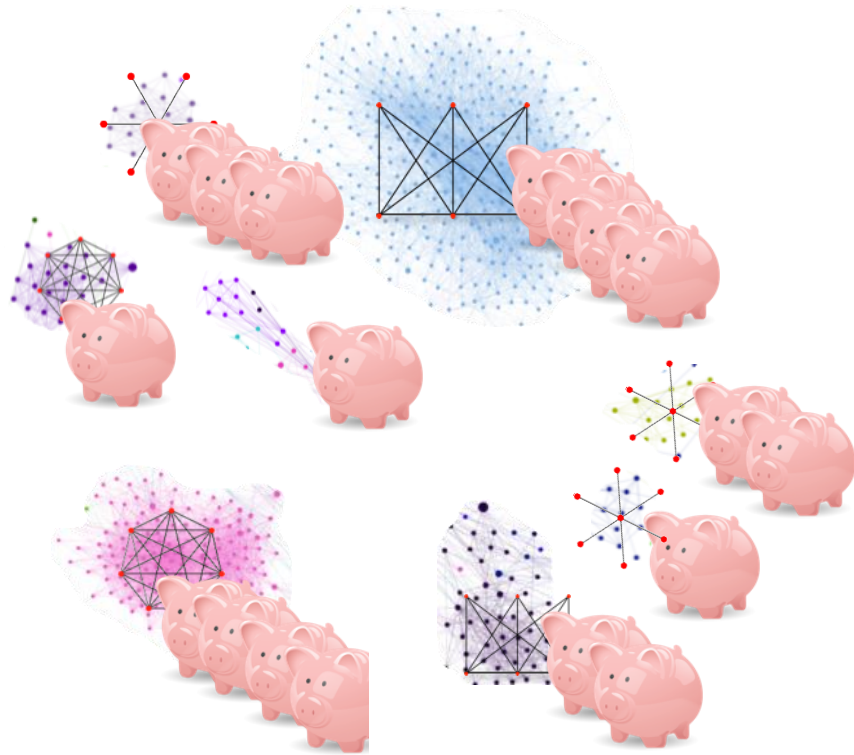
Concepts

Savings = # bits as structure - # bits as noise
compression gain



Step 3: Summary Assembly

TOP-k



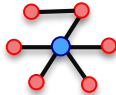


Summary

Concepts

Summary Encoding cost


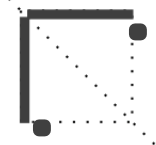
$$L(M) = \begin{matrix} \# \text{ of} \\ \text{structures} \end{matrix} + \begin{matrix} \# \text{ of} \\ \text{structures} \\ \text{per type} \end{matrix} + \begin{matrix} \text{for each structure} \\ \text{its encoding length} \end{matrix}$$

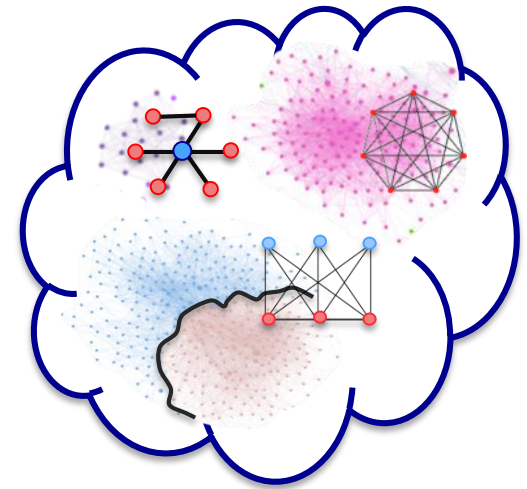
3

 : 1
 : 1
 : 1

 its
type

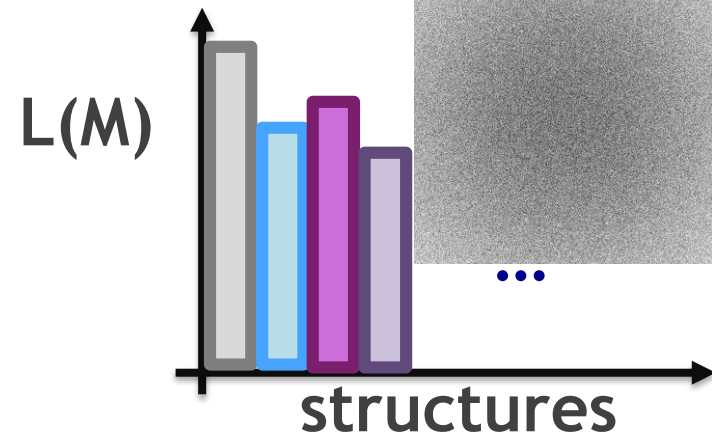
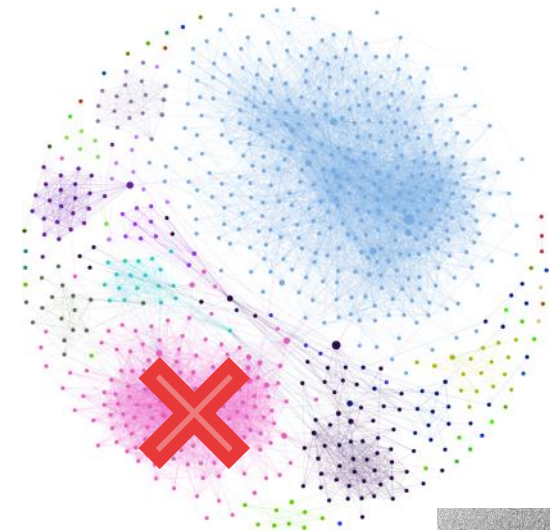
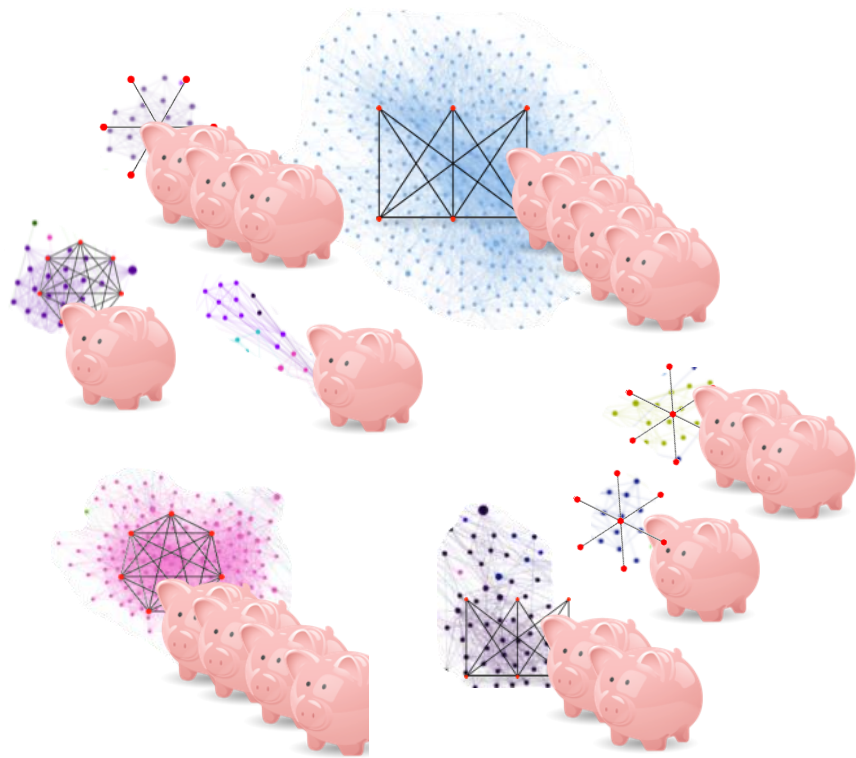
 its
connectivity



Step 3: Summary Assembly

Greedy&Forget



Roadmap

Main Idea

Encoding Schema

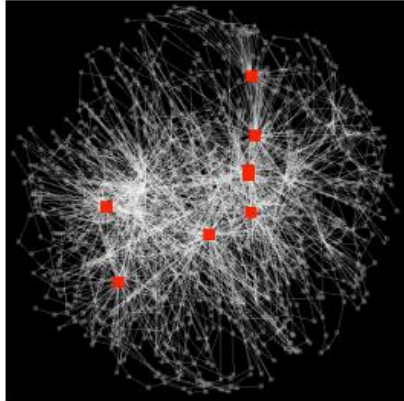
Proposed Algorithm: VoG

Experiments

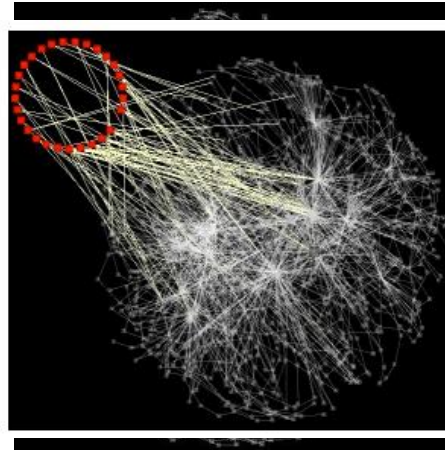
Conclusions



Application: Wikipedia controversy

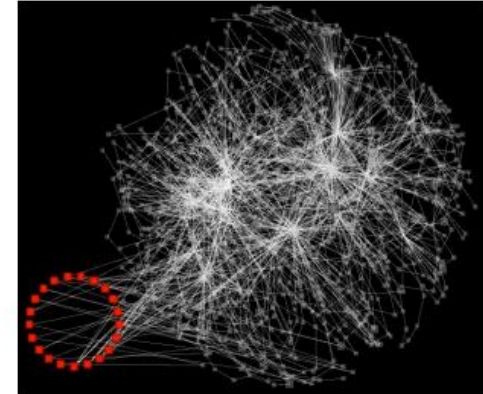


Stars:
admins,
bots,
heavy users



Bipartite cores: edit wars

↑
Kiev vs. Kyiv



↑
vandals

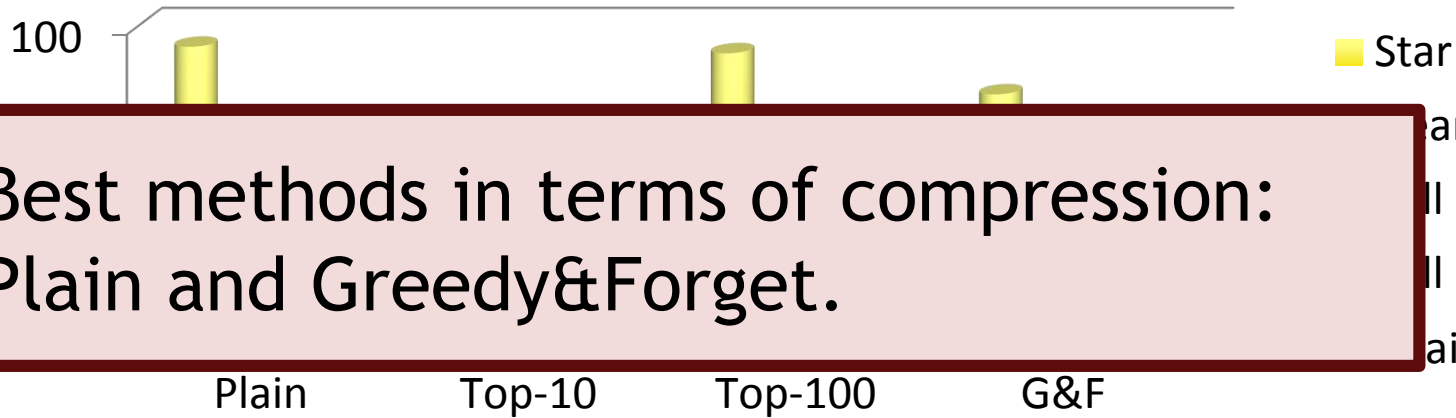
Nodes: wiki editors
Edges: co-edited



Quantitative Analysis



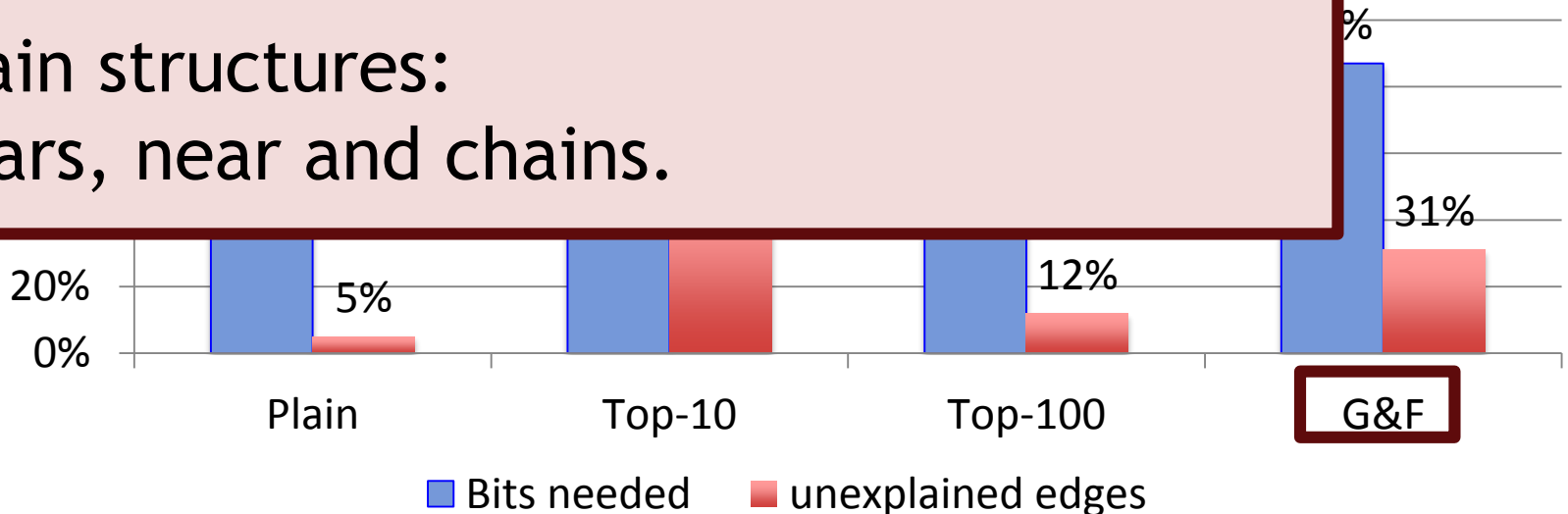
WIKIPEDIA
The Free Encyclopedia



Best methods in terms of compression:
Plain and Greedy&Forget.

19 833 bits as noise

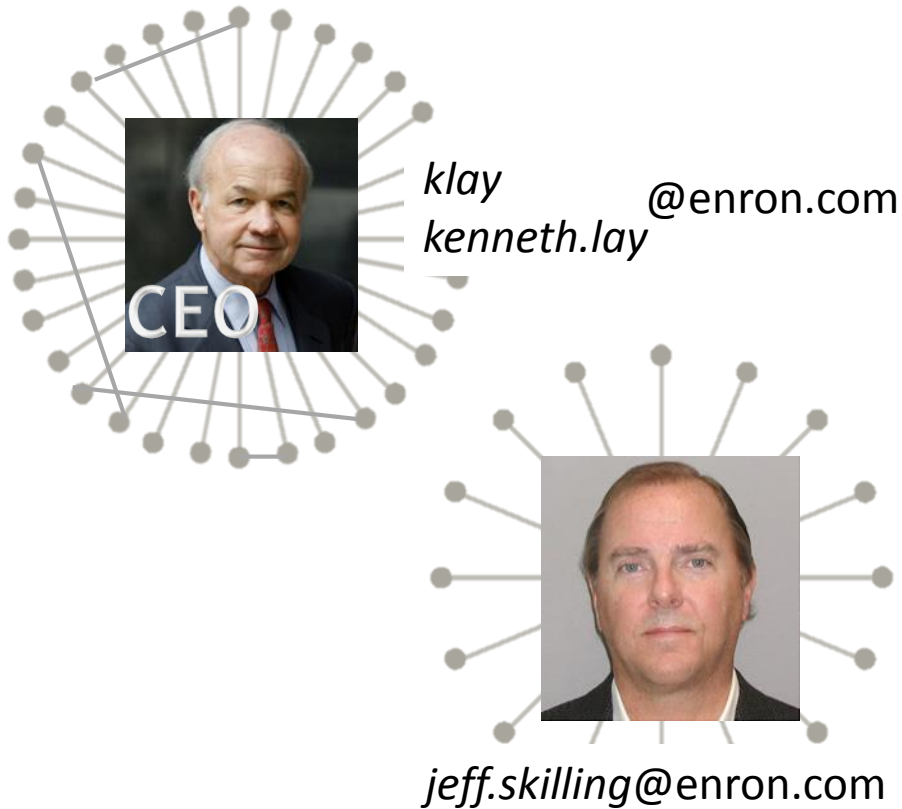
Main structures:
Stars, near and chains.



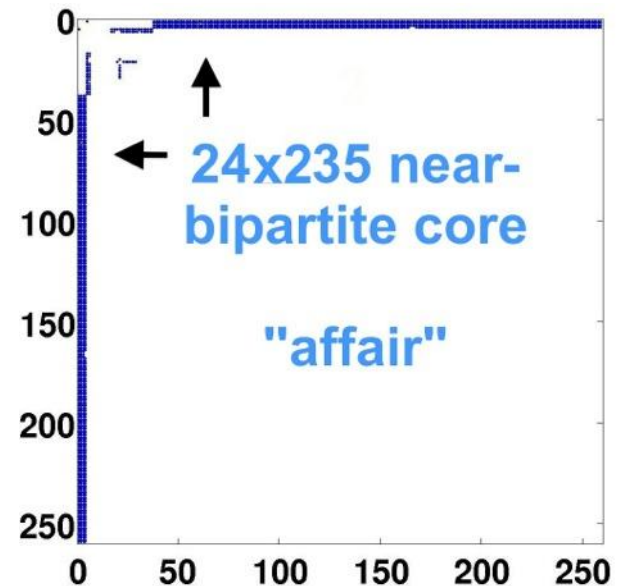
Application: Enron



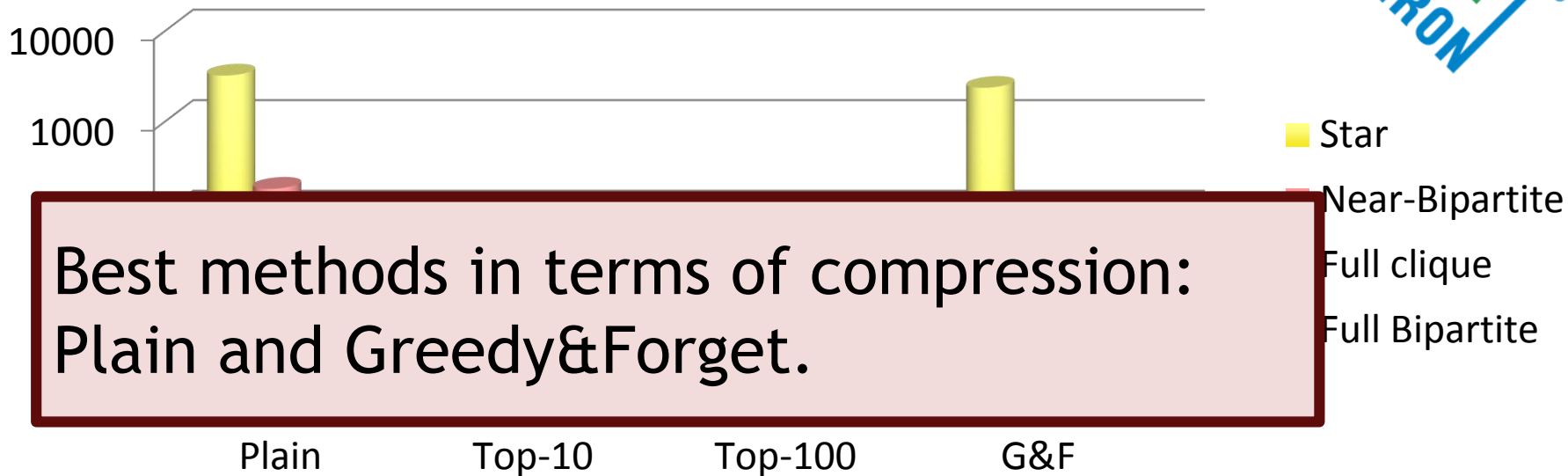
Top-3 Stars



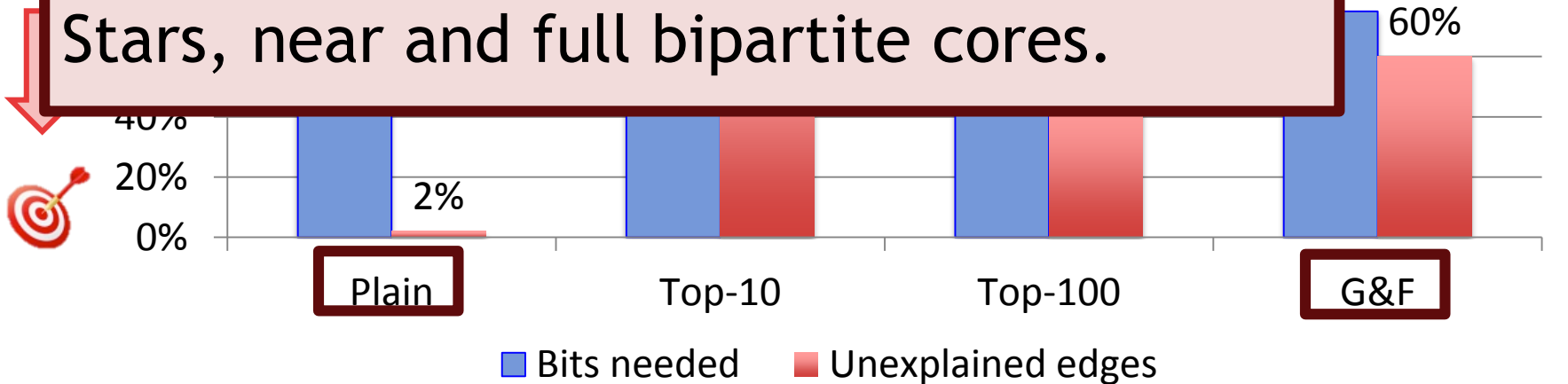
Top-1 NBC



Quantitative Analysis



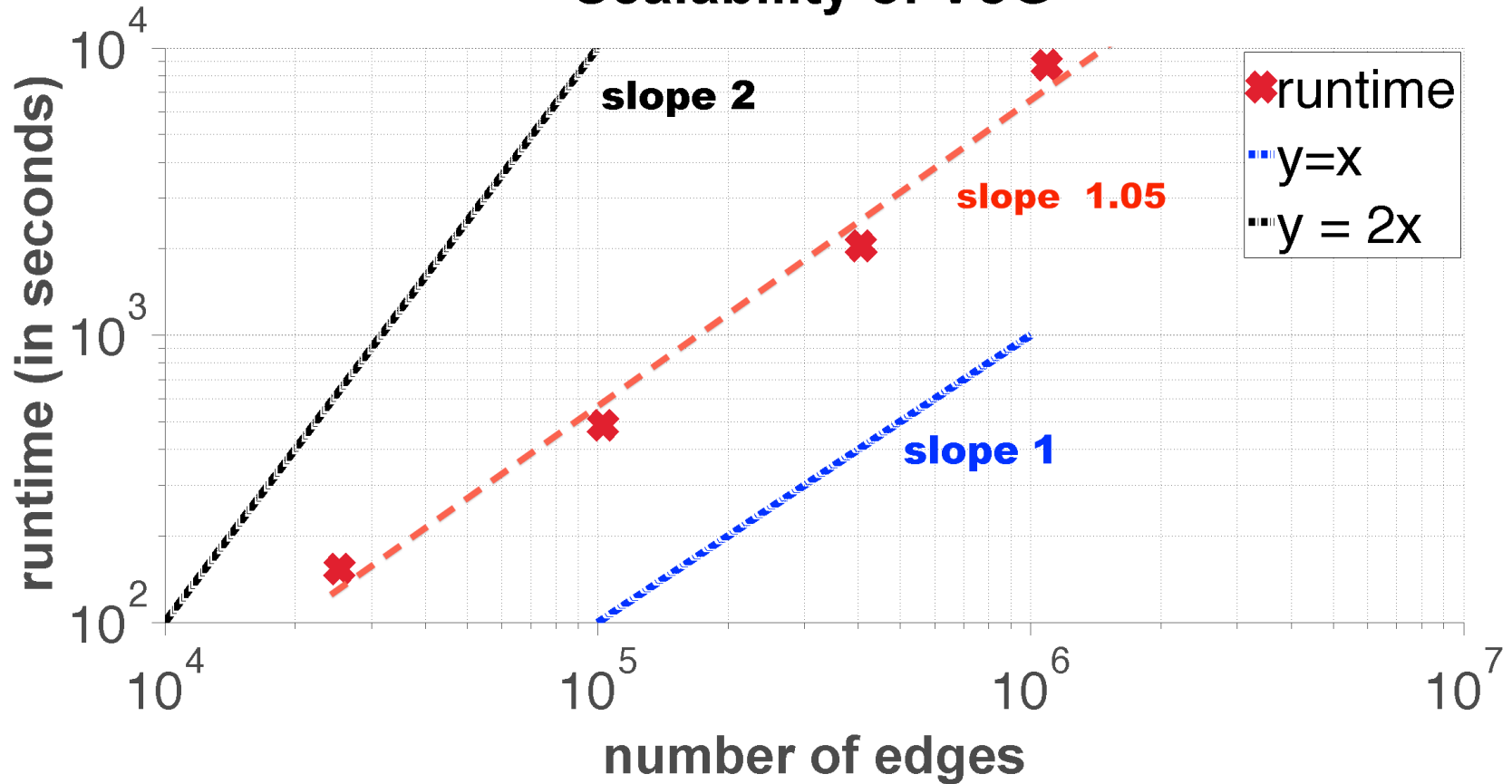
Main structures:
Stars, near and full bipartite cores.



Runtime



Scalability of VoG



VOG is *near-linear* on the number of edges of the input graph.



Roadmap

Main Idea

Encoding Schema

Proposed Algorithm: VoG

Experiments

Conclusions



Conclusions



- **Formulation:**
info-theoretic graph summarization approach
- **Algorithm:**
VoG is near-linear on the edges
- **Experiments on real graphs**

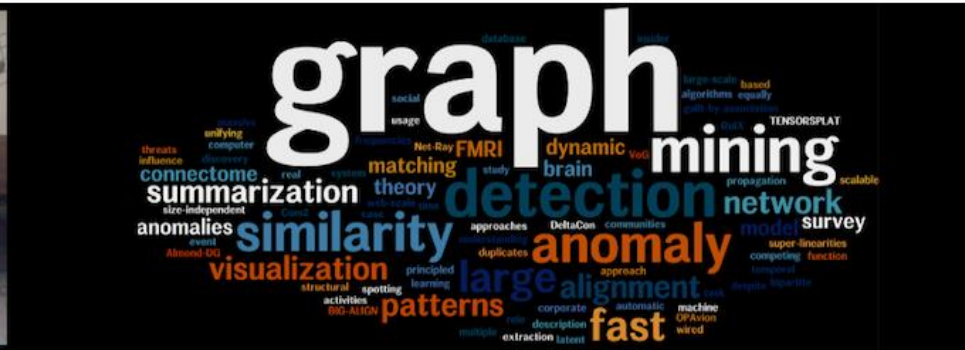


Code

www.cs.cmu.edu/~dkoutra/SRC/vog.tar




Danai Koutra



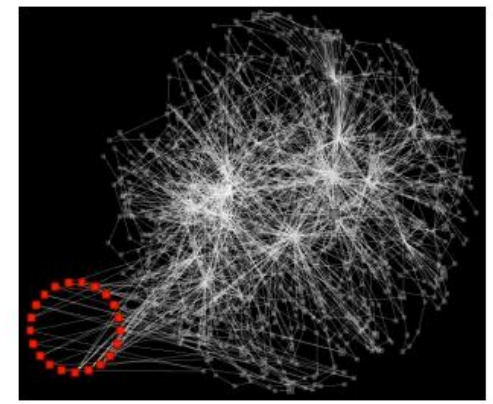
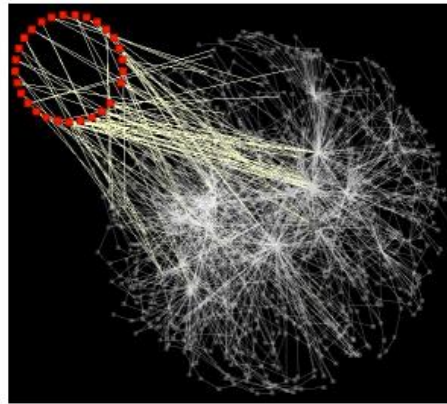
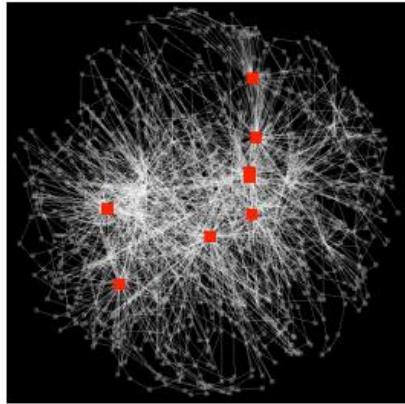
- About Me
- CV (pdf) - Jan. '14
- Bio
- Publications
- [Code New!](#)
- Contact
- Links
- Gallery

Code

-  **VoG: Summarizing and Understanding Large Graphs.**
Input: a graph
Output: a set of possibly overlapping subgraphs that most succinctly describe the given graph, i.e., that explain as many of its edges in as simple possible terms.
Publication: [VoG: Summarizing and Understanding Large Graphs.](#)



Thank you! Questions?



www.cs.cmu.edu/~dkoutra/pub.htm

danai@cs.cmu.edu

