

# Parameter-free Identification of Cohesive Subgroups in Large Attributed Graphs

Leman Akoglu



Carnegie Mellon University

Hanghang Tong



IBM T. J. Watson

Brendan Meeder



Carnegie Mellon University

Christos Faloutsos



Carnegie Mellon University

---

# PICS: problem

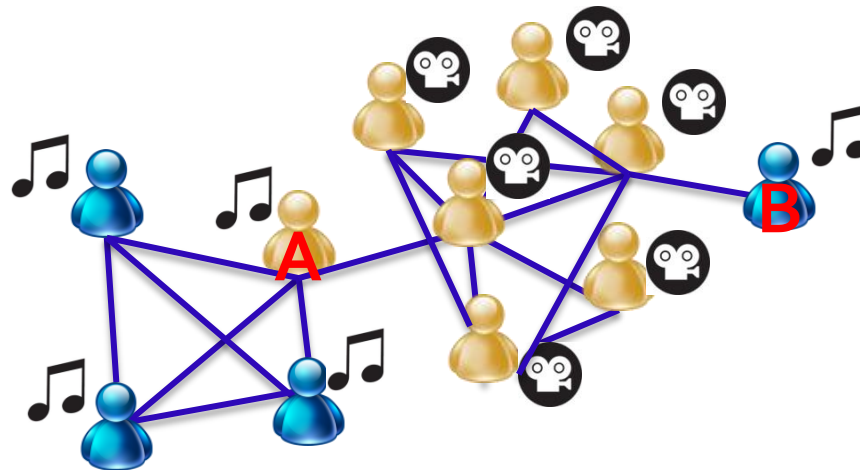
**Given** a graph with node attributes (features)

social networks + user interests

phone call networks + customer demographics

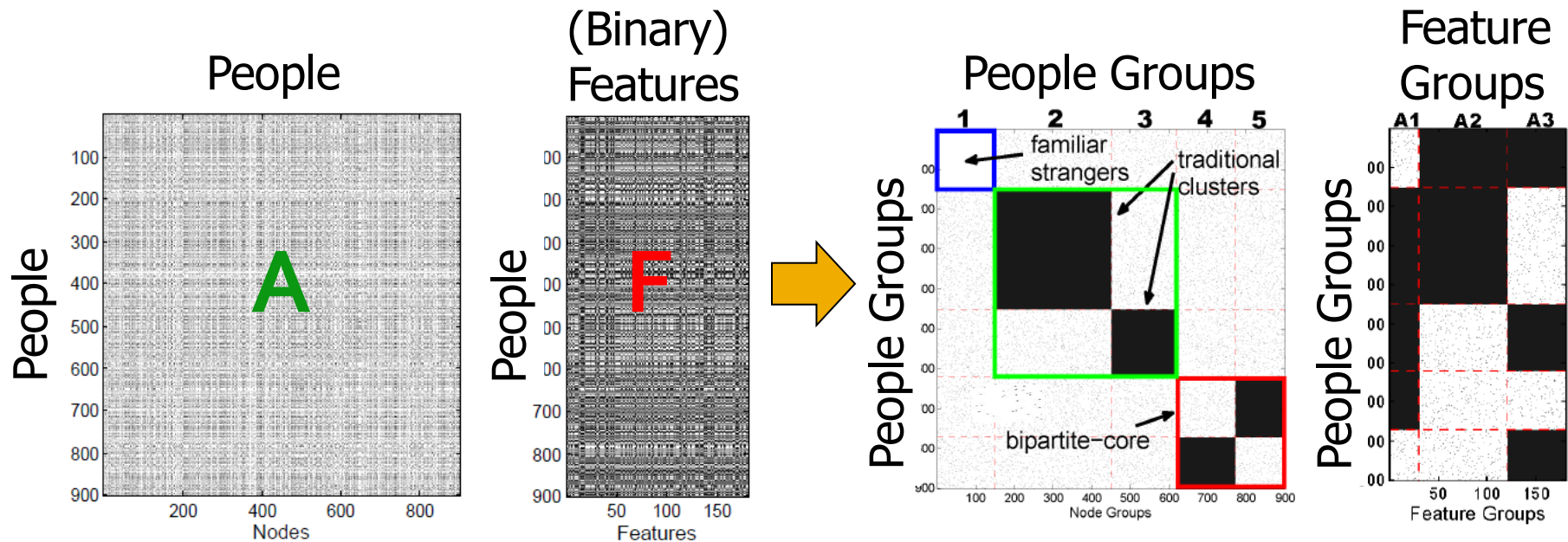
gene interaction networks + gene expression info

**Find** cohesive clusters, bridges, anomalies



**cohesive** cluster: similar connectivity & attribute coherence

# PICS: problem sketch

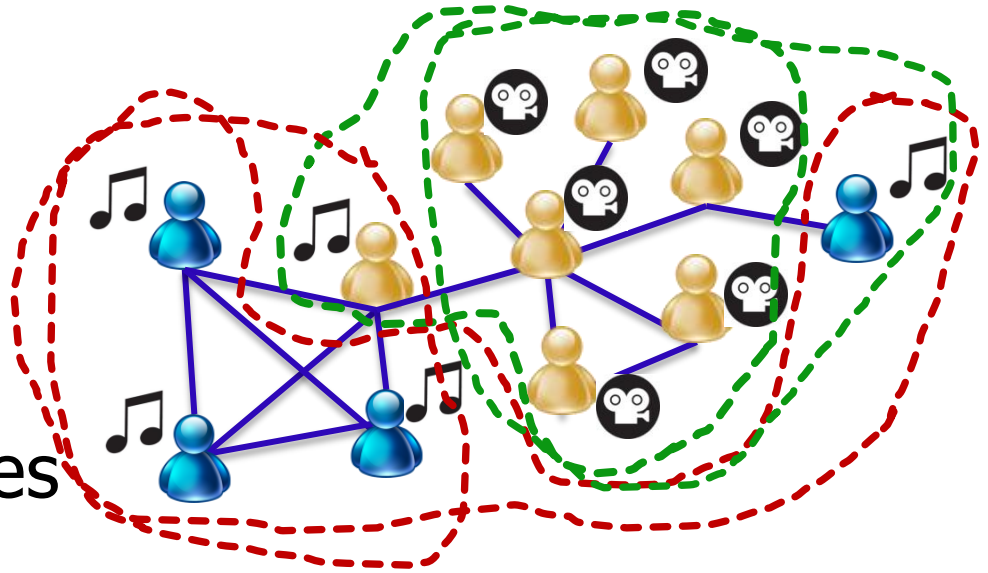


**Given** adjacency matrix  $A$  and feature matrix  $F$   
**Find** homogeneous blocks (clusters) in  $A$  and  $F$

- \* parameter-free
- \* scalable

# Simple extensions: why not?

- Flat clustering
- Graph clustering
- Additional feature nodes
  - heterogeneous graph
- Weighted edges by both connectivity and feature similarity
  - quadratic pairwise computations!
  - choice of similarity function



# Related Work

	Graph structure	Node attributes	Parameter-free	Linear scalability
Flat clustering (e.g. k-means) [Kriegel+] [Leeuwen+]		✓		✓
METIS [Karypis and Kumar], [Flake+] [Girvan and Newman] [Andersen+] spectral [Ng+], co-clustering [Dhillon+]	✓			✓
SA-cluster [Zhou+], Spect. rel. clus. [Long+]	✓	✓		
CoPaM [Moser+], Gamer [Gunneman+]	✓	✓		? , ✓
Autopart and cross-assoc.s [Chakrabarti+], GraphScope [Sun+], PaCK [He+]	✓		✓	✓

# PICS: approach

1. How many node- & attribute-clusters?
2. How to assign nodes and attributes to clusters?

Main idea: employ Minimum Description Length

$$\underbrace{L(M)} + \underbrace{L(D|M)}$$

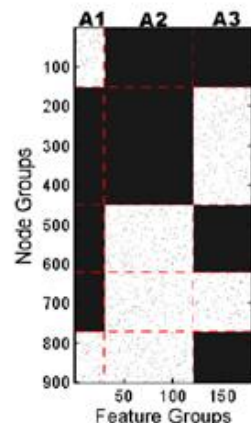
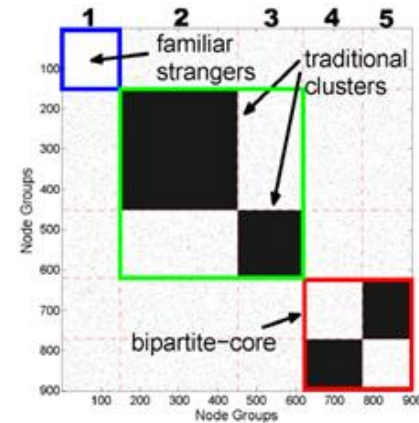
encoding length  
of clustering

encoding length  
of blocks

Good  
Clustering

implies

Good  
Compression



# Minimum Description Length

Given database  $D$  and set of models for  $D$ , MDL selects model  $M$  that minimizes

$$\underbrace{L(M)} + \underbrace{L(D|M)}$$

length in bits:  
description of  
model  $M$

length in bits: **data**,  
encoded by  $M$

↓

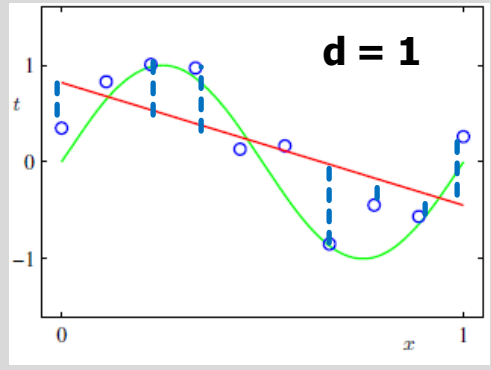
$$a_1x + a_0$$

↓

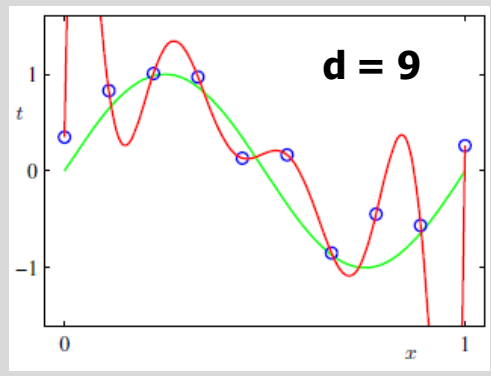
deltas

VS.

$$a_9x^9 + \dots + a_1x + a_0 \quad \{ \}$$



VS.



Bishop: PR&ML

# PICS: formulation

## ■ L (M) : Model description cost

1.  $\log^* n + \log^* f$     n: #nodes    f: #attributes

2.  $\log^* k + \log^* l$     k: #node-clus.    l: #attribute-clus.  
 $\log^*(k) = \log(k) + \log \log(k) + \dots$

3.  $nH(P) + fH(Q)$      $p_i = \frac{r_i}{n}$  ← size of node cluster  $i$   
 $q_j = \frac{c_j}{f}$  ← size of attr. cluster  $j$

$$\text{optimal \#bits} = -\log \frac{r_i}{n} = -\log p_i$$

$$\text{node clus. cost} = \sum_i r_i \cdot -\log \frac{r_i}{n} = n \cdot -\sum_i \frac{r_i}{n} \log \frac{r_i}{n} = nH(P)$$



# PICS: formulation

- $L(D|M)$ : Data description cost given Model


1. For each block in A and F , #1s:  $\log^* n_1(B_{ij})$

2. Encoding cost of a block

$$\begin{aligned} E(B_{ij}) &= -n_1(B_{ij}) \log_2(P_{ij}(1)) - n_0(B_{ij}) \log_2(P_{ij}(0)) \\ &= n(B_{ij}) H(P_{ij}(1)). \end{aligned}$$

where

$$P_{ij}(1) = n_1(B_{ij}) / n(B_{ij})$$


  
 $r_i C_j$  or  $r_i r_j$

# PICS: total cost objective

- **L (M) : Model description cost**

1.  $\log^* n + \log^* f$      $n$ : #nodes,  $f$ : #attributes
2.  $\log^* k + \log^* l$      $k$ : #node-clusters,  $l$ : #attribute-clusters
3.  $nH(P) + fH(Q)$      $p_i = \frac{r_i}{n}$  ← size of node-cluster  $i$   
 $c_j$  ← size of attribute-cluster  $j$

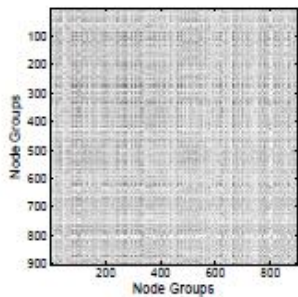
A similar problem (column re-ordering for minimum total run length) is shown to be NP-hard

- **[Johnson+]. (reduction from Hamiltonian Path)**

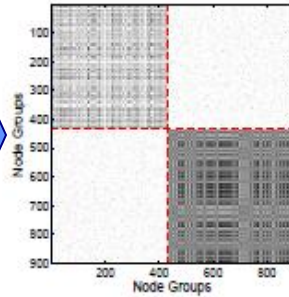
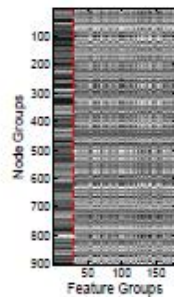
1.  $E(B_{ij}) = -n_1(B_{ij}) \log_2(P_{ij}(1)) - n_0(B_{ij}) \log_2(P_{ij}(0))$   
 $= n(B_{ij})H(P_{ij}(1)).$

where  $P_{ij}(1) = n_1(B_{ij})/n(B_{ij})$   
 $r_i c_j$  or  $r_i r_j$

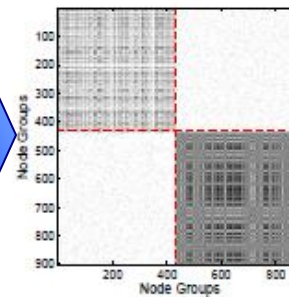
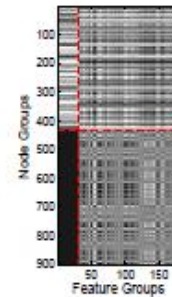
# PICS: algorithm sketch



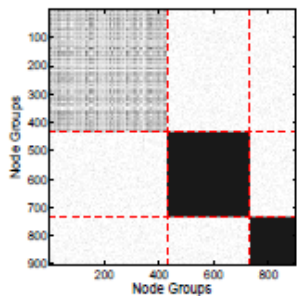
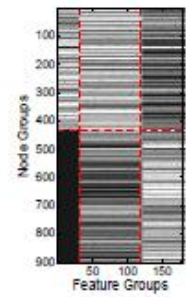
(a)  $k=1, l=2$   
Split-FeatureGroup



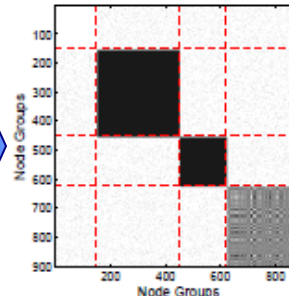
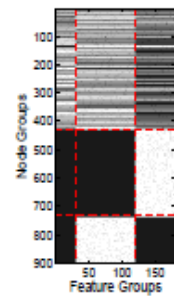
(b)  $k=2, l=2$   
Split-NodeGroup



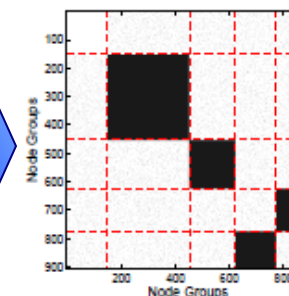
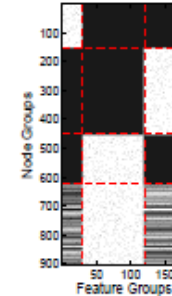
(c)  $k=2, l=3$   
Split-FeatureGroup



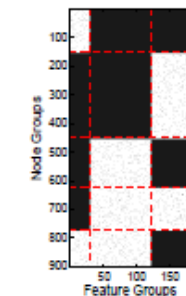
(d)  $k=3, l=3$   
Split-NodeGroup



(e)  $k=4, l=3$   
Split-NodeGroup



(f)  $k=5, l=3$   
Split-NodeGroup



The algorithm is iterative and monotonic  
—will converge to local optimum

# PICS: objective and algorithm

Total Encoding Cost (Length in bits)

$$L(\mathbf{A}, \mathbf{F}; R, C) = \log^* n + \log^* f + \log^* k + \log^* l$$

$$\begin{aligned} & - \sum_{i=1}^k r_i \log_2\left(\frac{r_i}{n}\right) - \sum_{j=1}^l c_j \log_2\left(\frac{c_j}{f}\right) \\ & + \sum_{i=1}^k \sum_{j=1}^l \left( \log^* n_1(B_{ij}^F) + E(B_{ij}^F) \right) \\ & + \sum_{i=1}^k \sum_{j=1}^k \left( \log^* n_1(B_{ij}^A) + E(B_{ij}^A) \right). \end{aligned}$$

Algorithm PICS

**Input:**  $n \times n$  link matrix  $\mathbf{A}$ ,  $n \times f$  feature matrix  $\mathbf{F}$

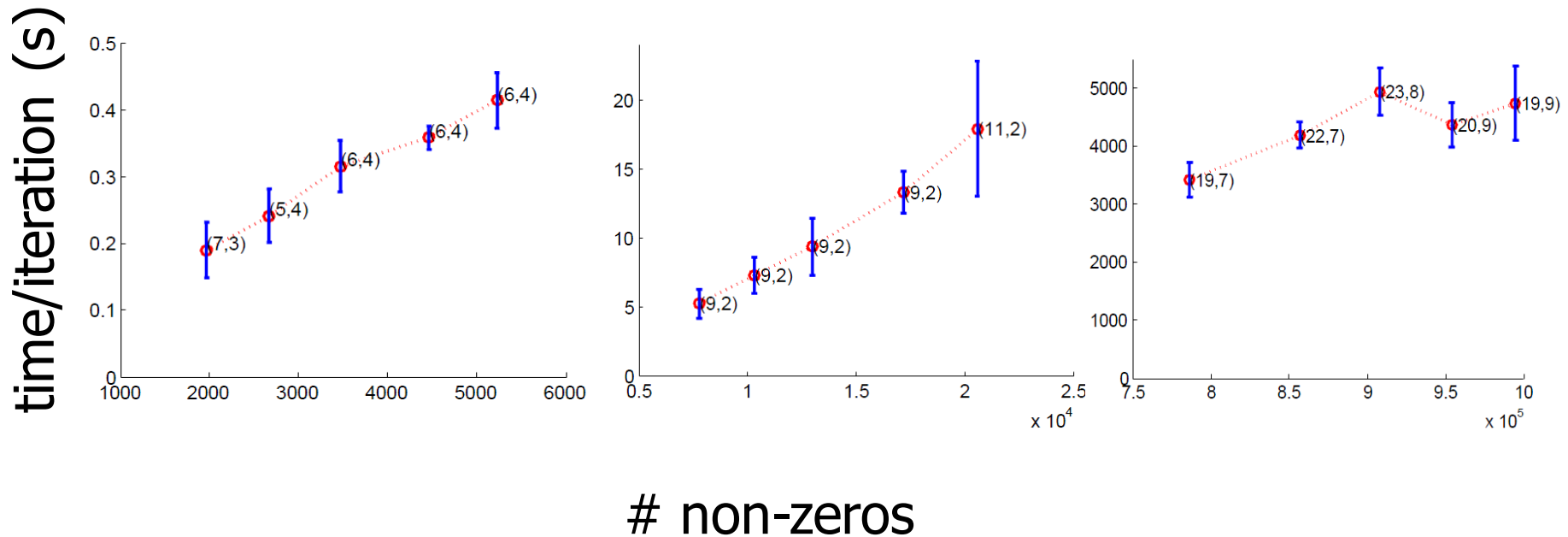
**Output:** A heuristic solution towards minimizing total encoding  $L(\mathbf{A}, \mathbf{F}; R, C)$ : number of row and column groups  $(k^*, l^*)$ , associated mapping  $(R^*, C^*)$

- 1: Set  $k^0 = l^0 = 1$  as we start with a single node and feature cluster.
- 2: Set  $R^0 := \{1, 2, \dots, n\} \rightarrow \{1, 1, \dots, 1\}$
- 3: Set  $C^0 := \{1, 2, \dots, f\} \rightarrow \{1, 1, \dots, 1\}$
- 4: Let  $T$  denote the outer iteration index. Set  $T = 0$ .
- 5: **repeat**
- 6:  $C^{T+1}, l^{T+1} := \text{Split-FeatureGroup}(\mathbf{F}, C^T, l^T)$
- 7:  $(R^{T+1}, C^{T+1}) := \text{Shuffle}(\mathbf{A}, \mathbf{F}, (R^T, C^{T+1}), (k^T, l^{T+1}))$
- 8:  $R^{T+1}, k^{T+1} := \text{Split-NodeGroup}(\mathbf{A}, \mathbf{F}, (R^{T+1}, C^{T+1}), (k^T, l^{T+1}))$
- 9:  $(R^{T+1}, C^{T+1}) := \text{Shuffle}(\mathbf{A}, \mathbf{F}, (R^{T+1}, C^{T+1}), (k^{T+1}, l^{T+1}))$
- 10: **if**  $L(\mathbf{A}, \mathbf{F}; R^{T+1}, C^{T+1}) \geq L(\mathbf{A}, \mathbf{F}; R^T, C^T)$  **then**
- 11:     **return**  $(k^*, l^*) = (k^T, l^T)$ ,  $(R^*, C^*) = (R^T, C^T)$
- 12: **else**
- 13:     Set  $T = T + 1$
- 14: **end if**
- 15: **until convergence**

# PICS: scalability

Computational complexity:

$$O(\max(k^*, l^*) * [2n_1(A)k^* + n_1(F)(k^* + l^*)] * \hat{t})$$



# PICS: datasets

## Graphs

1. **Phone call**

2. **Device**

3. **PolBooks**

4. **PolBlogs**

5. **Twitter**

6. **YouTube**

7. **YeastGene**

## Description

users, titles

users, titles

books, incl.

blogs, incl.

users, h-tags

users, groups

genes, articles

**n**

94

94

92

1.5K

9.6K

77K

844

**f**

7

7

2

2

10K

30K

17K

**nnz**

391

5K

840

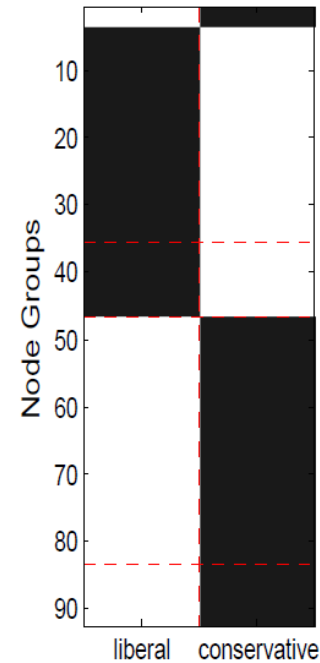
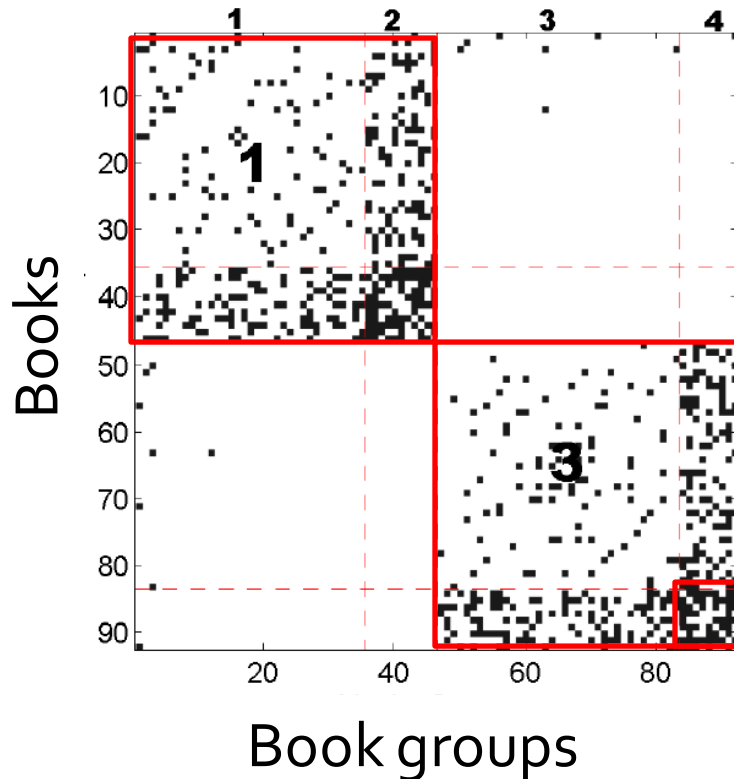
20K

82K

1M

64K

# PICS at work (Political books)



liberal vs.  
conservative

“core and periphery”

# PICS at work (Political books)

## Examples of “core” liberal and conservative books



### Liberal

- *Lies and the Lying Liars Who Tell Them: A Fair and Balanced Look at the Right*
- *Big Lies: The Right-Wing Propaganda Machine and How It Distorts the Truth*
- *The Lies of George W. Bush*
- *Dude, Where's My Country?*

### Conservative

- *Persecution: How Liberals Are Waging War Against Christianity*
- *Deliver Us from Evil: Defeating Terrorism, Despotism, and Liberalism*
- *Tales from the Left Coast*
- *A National Party No More*

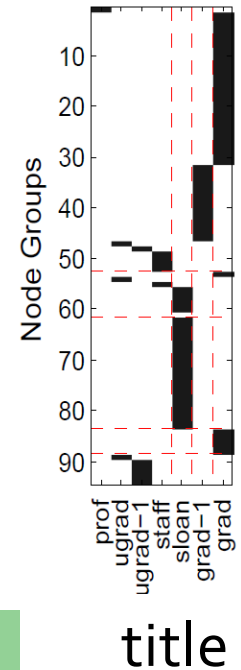
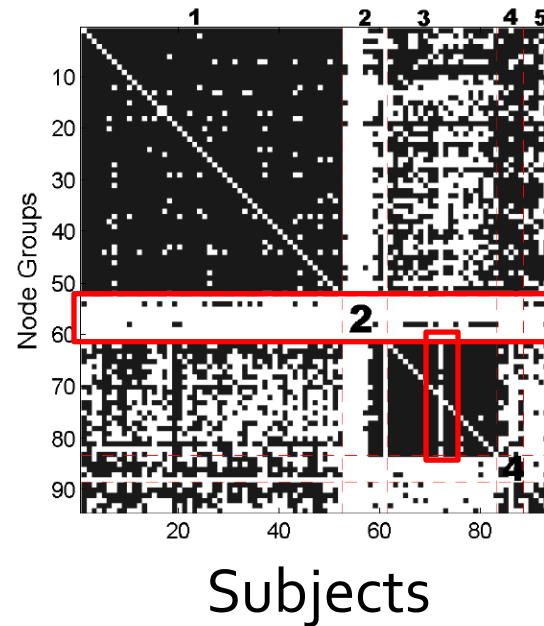
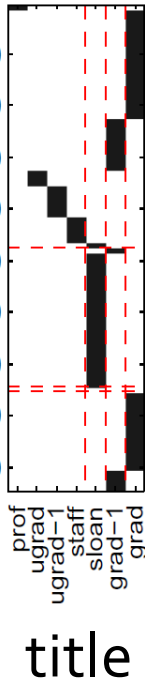
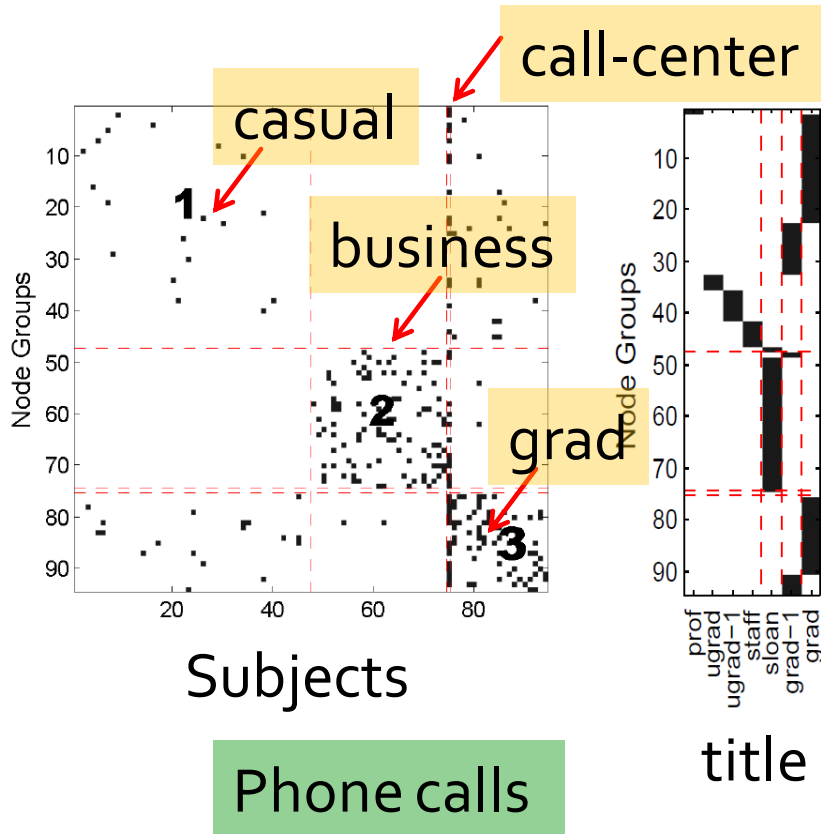
## Examples of bridging ‘conservative’ books

- *Bush at War*
- *The Bushes: Portrait of a Dynasty*
- *Rise of the Vulcans: The History of Bush's War Cabinet*

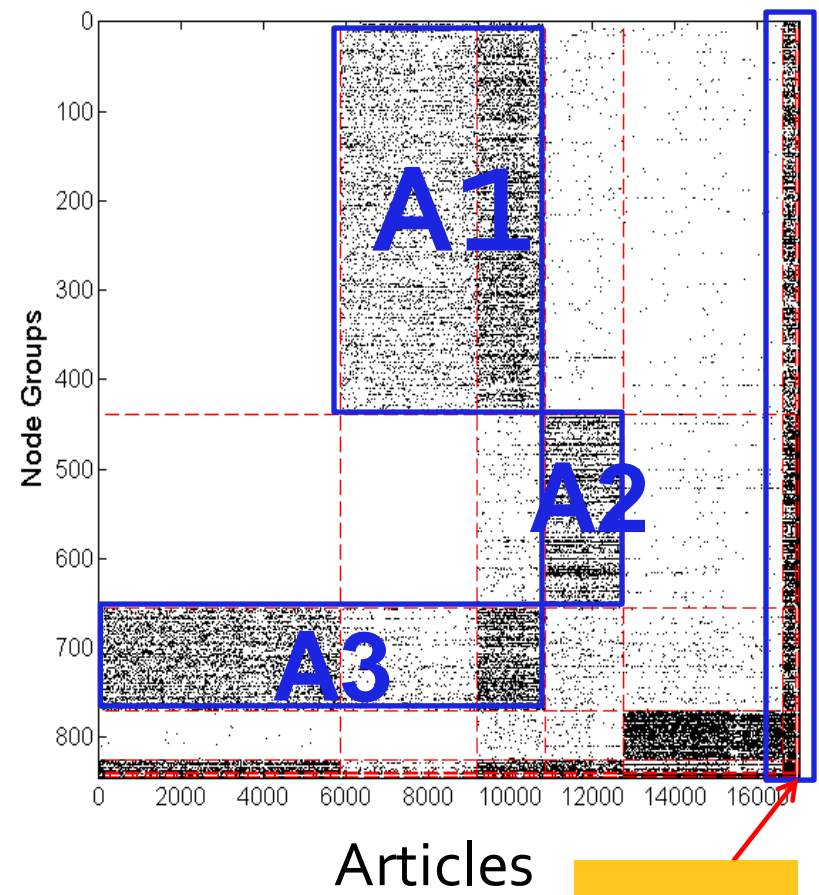
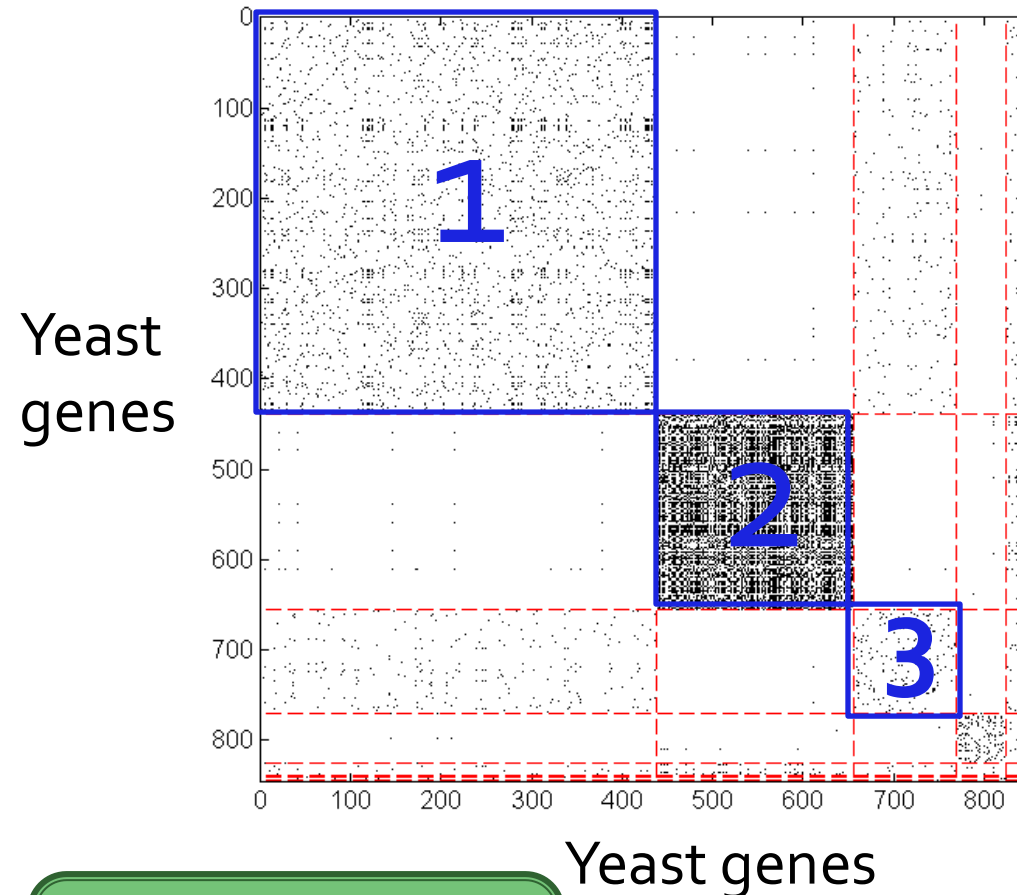
“core and periphery”



# PICS at work (Reality mining)

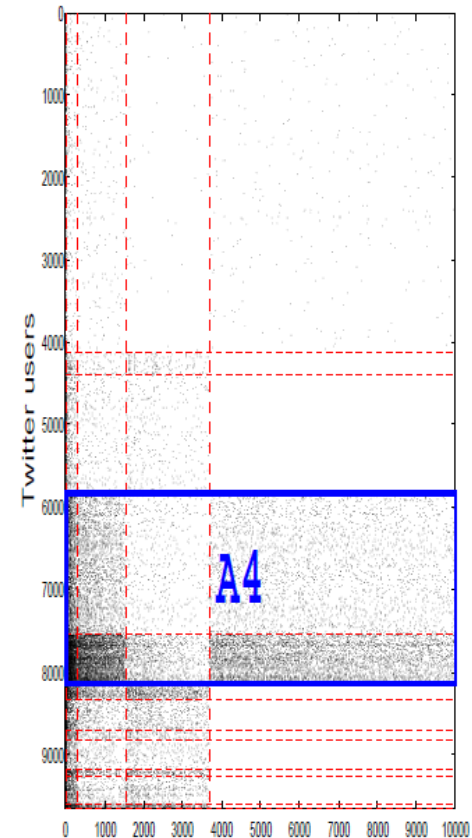
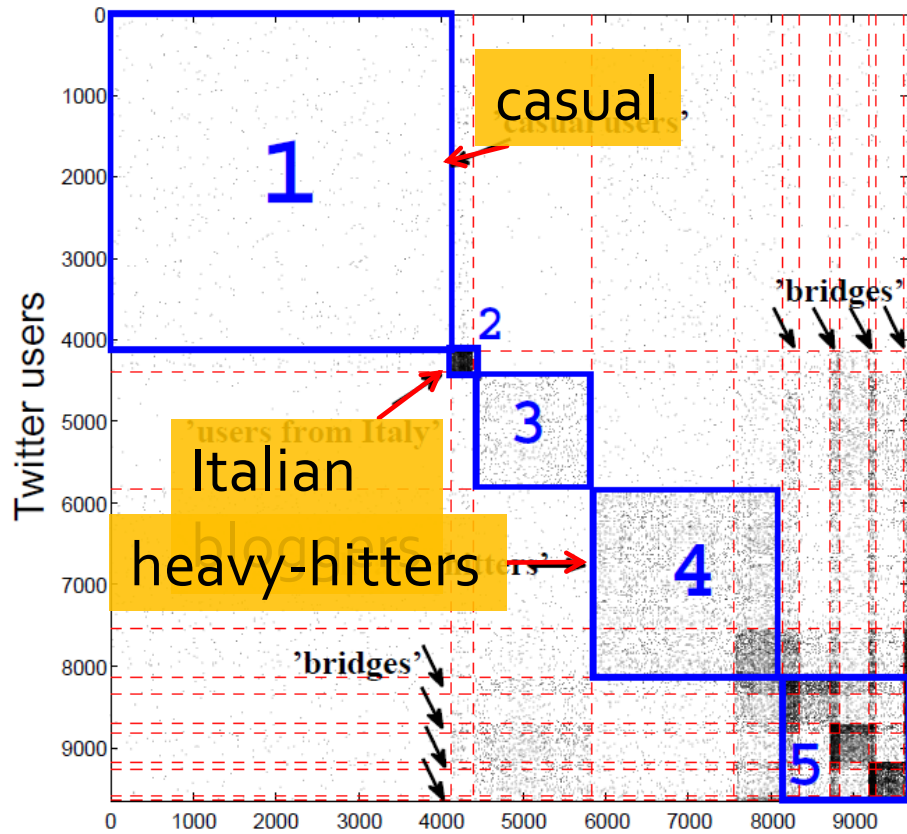


# PICS at work (YeastGene)



844 genes  
17K articles

# PICS at work (Twitter)



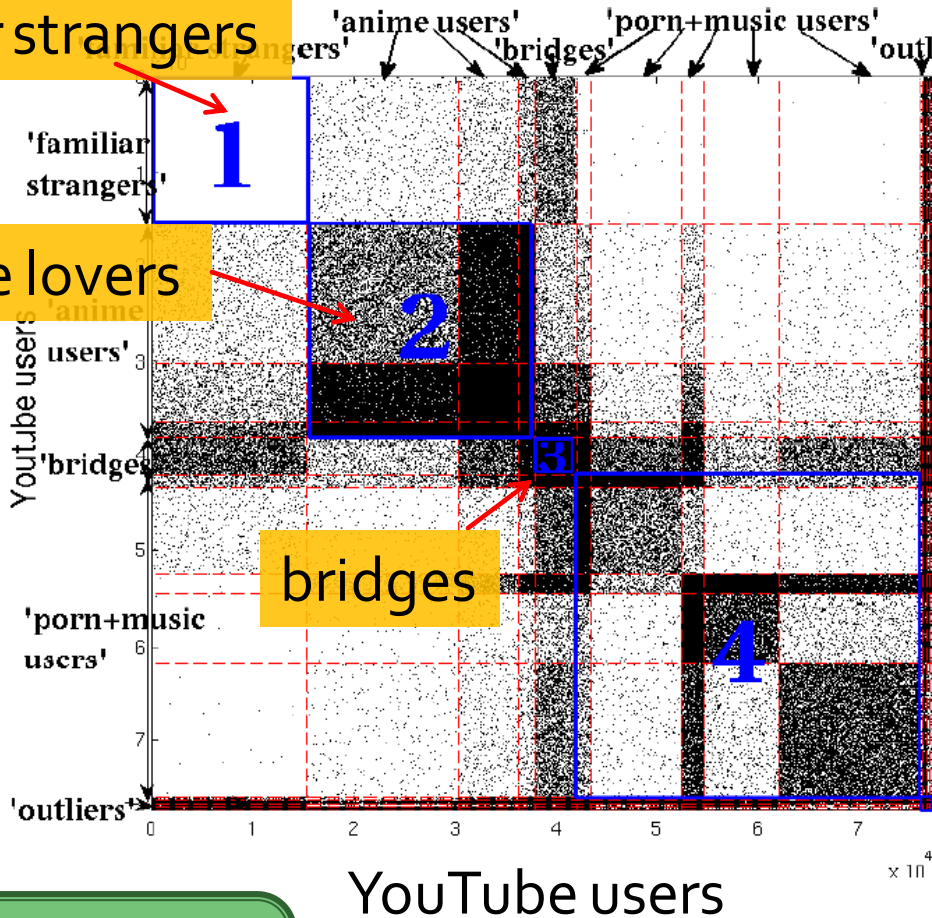
Twitter users

@hashtags

9,6K users  
10K hashtags

# PICS at work (YouTube)

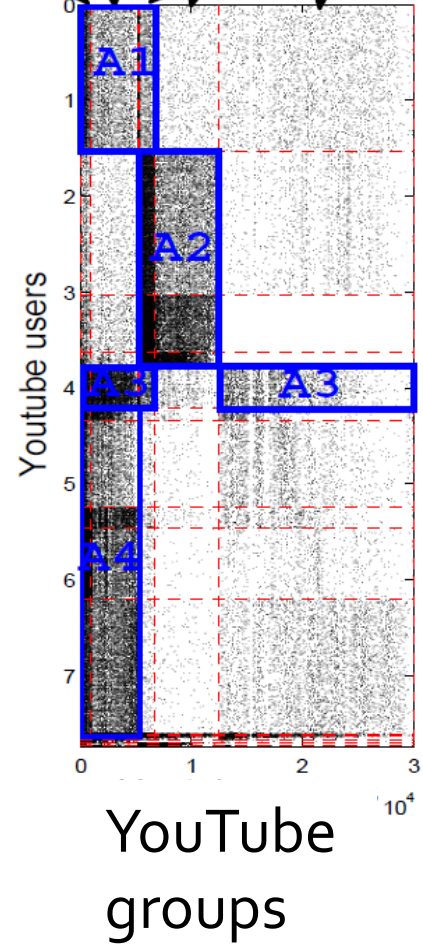
familiar strangers



anime lovers

bridges

porn music anime general interest



77K users  
30K groups

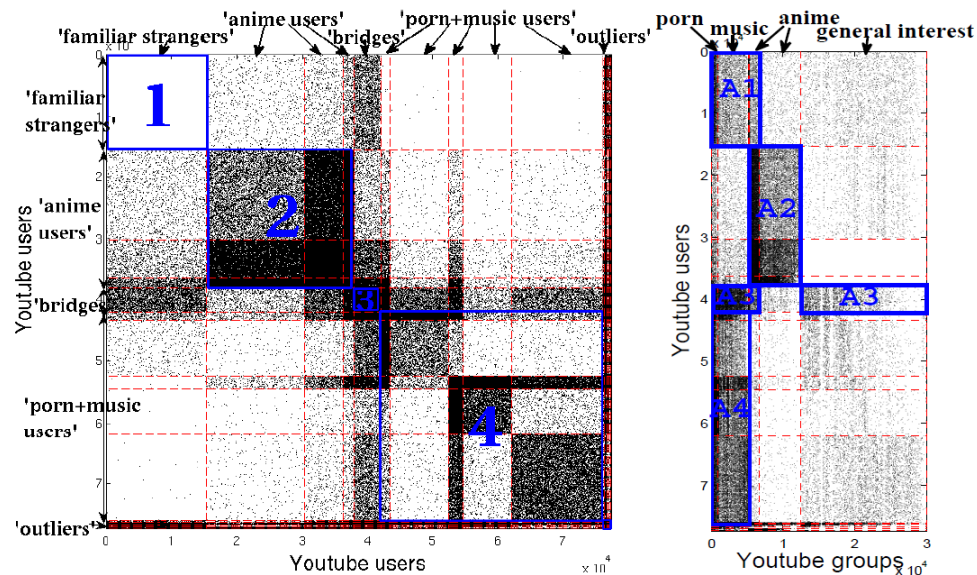
# Summary of contributions

- **Novel clustering model:**
  - PICS finds groups of nodes in an attributed graph with (1) similar connectivity, and (2) attribute homogeneity.
  - It also groups the node attributes into attribute-clusters.
- **Parameter-free nature:**
  - No user input, e.g. number of clusters, similarity functions/thresholds
- **Effectiveness:**
  - Insightful clusters, bridges and outliers in diverse real-world datasets including YouTube and Twitter.
- **Scalability:**
  - Linearly growing run time with graph + attribute size

# Thank you!

[lakoglu@cs.cmu.edu](mailto:lakoglu@cs.cmu.edu)

<http://www.cs.cmu.edu/~lakoglu/>



Source code: [www.cs.cmu.edu/~lakoglu/#pics](http://www.cs.cmu.edu/~lakoglu/#pics)