



# Social network analysis: Community detection

Rushed Kanawati

LIPN, CNRS UMR 7030; USPC

<http://lipn.fr/~kanawati>

[rushed.kanawati@lipn.univ-paris13.fr](mailto:rushed.kanawati@lipn.univ-paris13.fr)





# COMPLEX NETWORK

## Definition

Graphs modeling (direct/indirect) interactions among actors.

## Basic topological features

- ▶ Low Density
- ▶ Small Diameter
- ▶ Heterogeneous degree distribution.
- ▶ High Clustering coefficient
- ▶ Community structure





## EXAMPLE II: SPATIAL NETWORK

## Public sites accessibility in the Bourget district

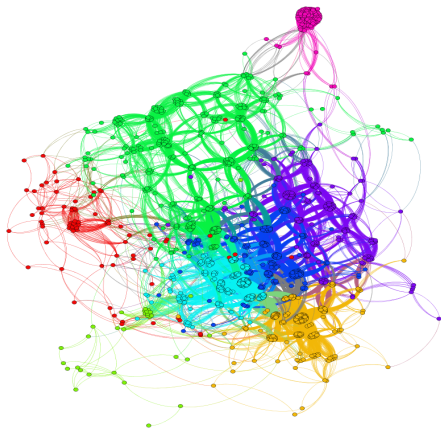
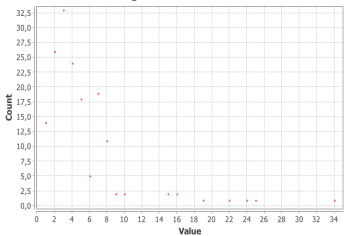
projet FUI UrbanD 2009-2012

Densiy : 0.052

Diameter : 7

clustering Coef.: 0.87

Degree Distribution



A relative-neighbourhood graph

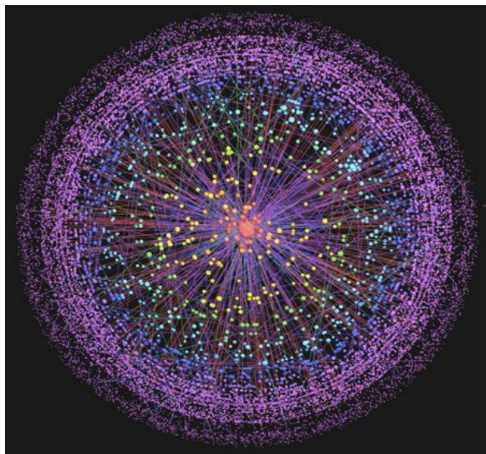






# ANOTHER EXAMPLE

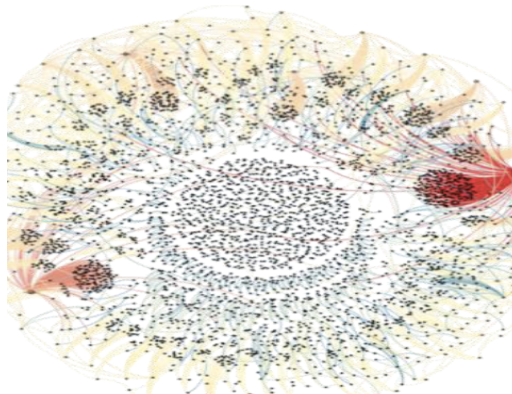
- ▶ Graph of the Internet
- ▶ Nodes = service providers
- ▶ Edges : Connections



[CHK<sup>+</sup>07]

# ANOTHER EXAMPLE

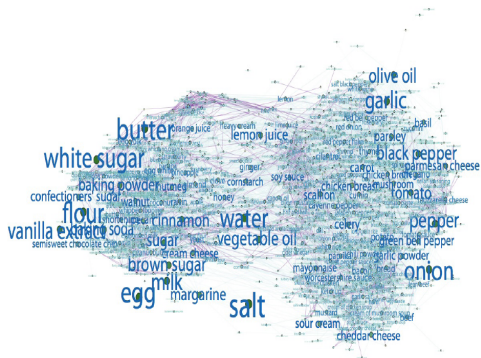
- ▶ Mubarak's resignation announcement.
- ▶ Nodes = Twitter user
- ▶ Edges : retweet on #jan25 hashtag



<http://gephi.org/2011/the-egyptian-revolution-on-twitter/>

# LAST EXAMPLE !

- ▶ Ingredients network extracted from cooking receipts
- ▶ Nodes = Ingredients
- ▶ Edges : usage in the same receipt.



[Lada Adamic Course on Network Analysis: Coursera]



# COMPLEX NETWORK ANALYSIS: TASKS

## Node oriented

- ▶ Nodes ranking in function of their importance, influence, ...
- ▶ Applying *centrality* functions
- ▶ Applications: Ranking (ex. researchers ! ), Viral Marketing, Cyber-attacks (prevention); ...

## Network oriented

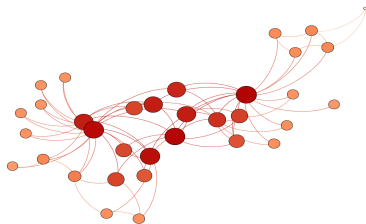
- ▶ Network formation & evolution models
- ▶ Link prediction, Diffusion models
- ▶ Applications: Explanation, Recommendation, Anomalies detection, information diffusion, ...

## Community oriented



# CENTRALITY: SOME EXAMPLES

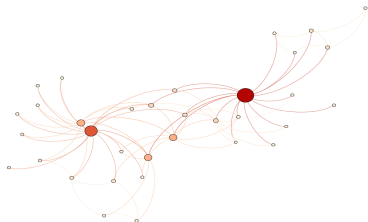
## Closeness centrality



- ▶  $C_c(v) = \frac{1}{\sum_{u \in V} sp(v,u)}$
- ▶  $sp(u, v)$  : shortest path length between  $u, v$ .
- ▶ Complexity =  
 $\mathcal{O}(n \times m + n^2 \log n)$

# CENTRALITY: SOME EXAMPLES

## Betweenness centrality



- ▶  $C_i(v) = \sum_{s,t \in V, stv} \frac{\sigma_{s,t}(v)}{\sigma_{s,t}}$
- ▶  $\sigma_{s,t}(v)$  : number of shortest paths linking  $s$  to  $t$  that include  $v$
- ▶  $\sigma_{s,t}$  : total number of shortest paths linking  $s$  to  $t$
- ▶ Complexity =  $\mathcal{O}(n^3)$

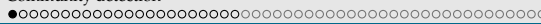




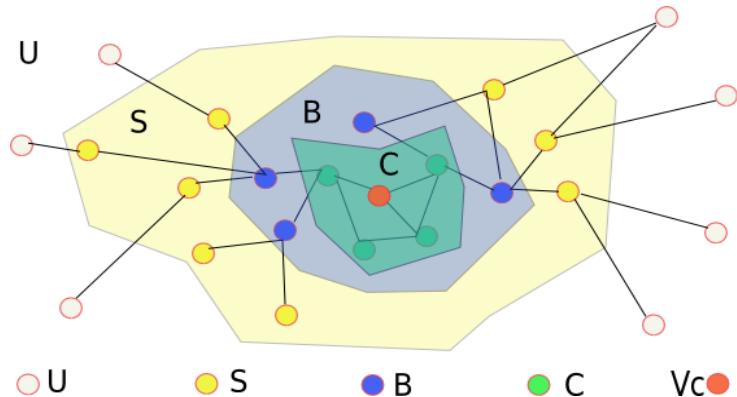
# COMMUNITY DETECTION PROBLEM

- ▶ Local community identification (ego-centred).
- ▶ Network partition computing
- ▶ *Overlapping community detection*





# LOCAL COMMUNITY





# LOCAL COMMUNITY

- 1  $C \leftarrow \{\phi\}, B \leftarrow \{n_0\} S \leftarrow \Gamma(n_0)$
- 2  $Q \leftarrow 0$  /\* a community **quality function** \*/
- 3 While  $Q$  can be enhanced Do
  - 1  $n \leftarrow \operatorname{argmax}_{n \in S} Q$
  - 2  $S \leftarrow S - \{n\}$
  - 3  $D \leftarrow D + \{n\}$
  - 4 update  $B, S, C$
- 4 Return  $D$





# MULTI-OBJECTIVE LOCAL COMMUNITY IDENTIFICATION [?]

## Three main approaches

Combine then Rank

**Ensemble ranking**

Ensemble clustering

















# ENSEMBLE RANKING : APPROACHES



John Kemeny 1926-1992

## Optimal Kemeny rank aggregation

- ▶ Let  $d()$  be distance over rankings  $\sigma_i$  (ex. Kendall  $\tau$ , Spearman's footrule)
- ▶ **Find  $\pi$  that minimise  $\sum_i d(\pi, \sigma_i)$**
- ▶ NP-Hard problem
- ▶ Approximation : Local Kemeny : two adjacent candidates are in the good order.
- ▶ **Local Kemeny** : Apply Bubble sort using the *majority preference partial order relationship*
- ▶ **Approximate Kemeny** : Apply QuickSort

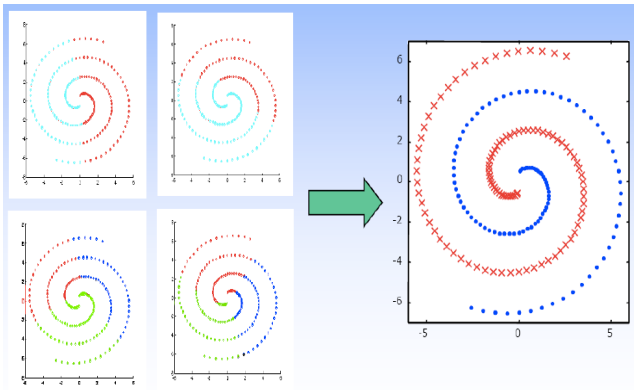


# ENSEMBLE CLUSTERING APPROACHES

## Principle

- ▶ Let  $C_{v_q}^{Q_i}$  be the the local community of  $v_q$  applying  $Q_i$ .
- ▶ We have a natural partition :  $P_{Q_i} = \{C_{v_q}^{Q_i}, \overline{C_{v_q}^{Q_i}}\}$
- ▶ **Apply an ensemble clustering approach.**

# ENSEMBLE CLUSTERING: PRINCIPLE



from A. Topchy et. al. Clustering Ensembles: Models of Consensus and Weak Partitions. PAMI, 2005





# ENSEMBLE CLUSTERING: APPROACHES

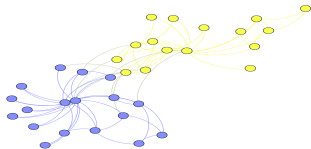
## HGPA: HyperGraph-Partitioning Algorithm

- ▶ Construct a hypergraph where nodes are objects and hyperedges are clusters.
- ▶ Partition the hypergraph by minimizing the number of cut hyperedges
- ▶ Each component forms a meta cluster
- ▶ Complexity :  $\mathcal{O}(nkr)$

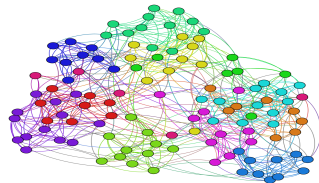




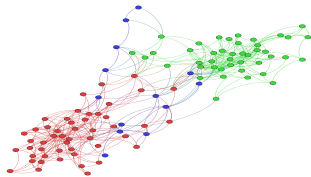
# DATASETS



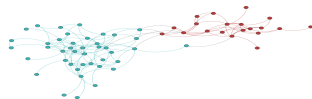
Zachary's Karate Club



Football network



US Political books network

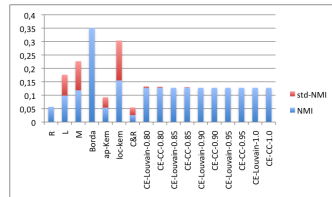
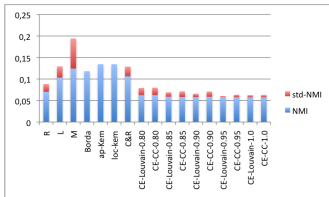


Dolphins social network

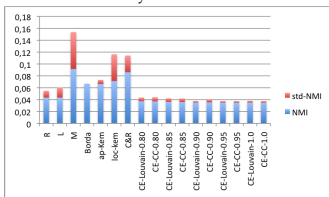




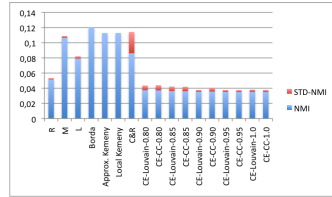
# RESULTS : COMPARATIVE RESULTS (NMI)



### Zachary's Karate Club



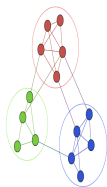
### Football network



### US Political books network

### Dolphins social network

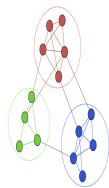
# GLOBAL COMMUNITIES DETECTION



## Problem

- ▶ Divid the set of nodes in a number of (overlapping) subsets such that induced subgraphs are dense and loosley coupled.

# GLOBAL COMMUNITIES DETECTION



## Problem

- ▶ Divid the set of nodes in a number of (overlapping) subsets such that induced subgraphs are dense and loosley coupled.

## Recommended readings

- ▶ **S. Fortunato.** *Community detection in graphs*. Physics Reports, 2010, 486, 75-174
- ▶ **L. Tang, H. Liu.** *Community Detection and Mining in Social Media*, Morgan Claypool Publishers, 2010
- ▶ **R. Kanawati,** *Détection de communautés dans les grands graphes d'interactions multiplexes : état de l'art*, (RNTI 2014), hal-00881668





# GROUP-BASED APPROCHES

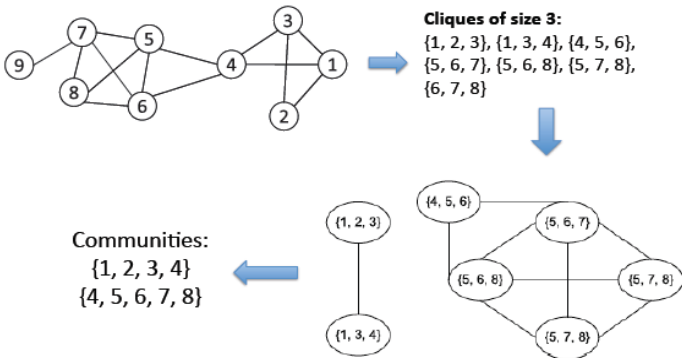
## Principle

Search for special (dense) subgraphs:

- ▶ k-clique
- ▶ n-clique
- ▶  $\gamma$ -dense clique
- ▶ K-core



# EXAMPLE: CLIQUE PERCOLATION



Suits fairly dense graph.

# NETWORK-BASED APPROCHES

## Clustering approaches

- ▶ Apply classical clustering approaches using *graph-based* distance function
- ▶ Different types of Graph-based distances: neighborhood-based, path-based (Random-walk)
- ▶ **Usually requires the number of clusters to discover**

# MODULARITY OPTIMIZATION APPROACHES

Modularity: a partition quality criteria

$$Q(\mathcal{P}) = \frac{1}{2m} \sum_{c \in \mathcal{P}} \sum_{i,j \in c} (A_{ij} - \frac{d_i d_j}{2m}) \quad (1)$$

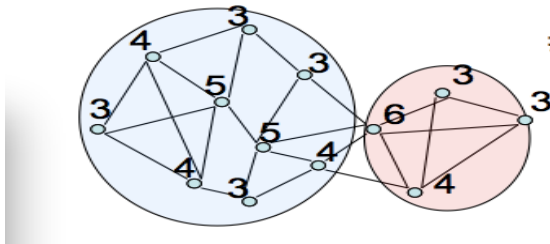


Figure: Example :  $Q = \frac{(15+6) - (11.25+2.56)}{25} = 0.275$

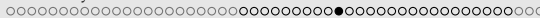
# MODULARITY OPTIMIZATION APPROACHES

- ▶ Applying classical optimization algorithms (ex. Genetic algorithms [Piz12]).
- ▶ Applying hierarchical clustering and select the level with  $Q_{max}$  (ex. Walktrap [PL06])
- ▶ Divisive approach : Girvan-Newman algorithm [GN02]
- ▶ Greedy optimization : Louvain algorithm [BGL08]
- ▶ ...

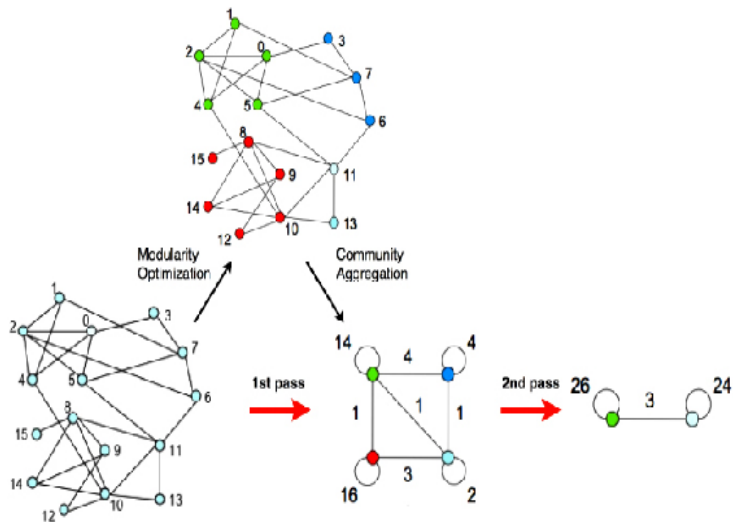
# EXAMPLE: GIRVAN-NEWMAN ALGORITHM

- 1 Compute betweenness centrality for each edge.
- 2 Remove edge with highest score.
- 3 Re-compute all scores.
- 4 Repeat 2nd step.

**Complexity** :  $\mathcal{O}(n^3)$



# EXAMPLE: LOUVAIN APPROACH



# MODULARITY OPTIMIZATION LIMITATIONS

## Hypothesis

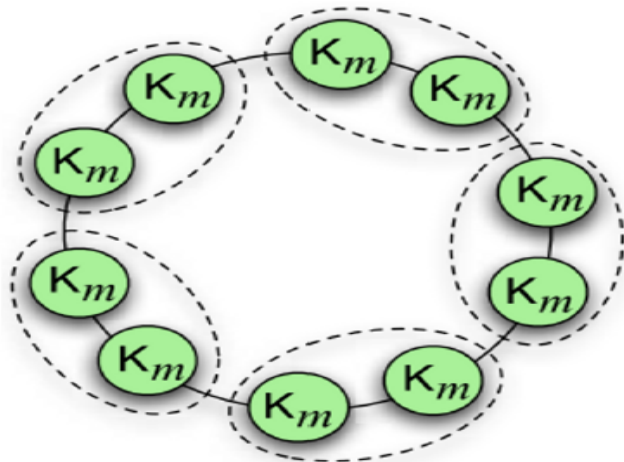
The best partition of a graph is the one that maximize the modularity.

If a network has a community structure, then it is popsicle to find a precise partition with maximal modularity

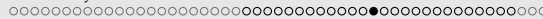
If a network has a community structure, then partitions inducing high modularity values are structurally similar.

All three hypothesis do not hold [GdMC10, LF11].

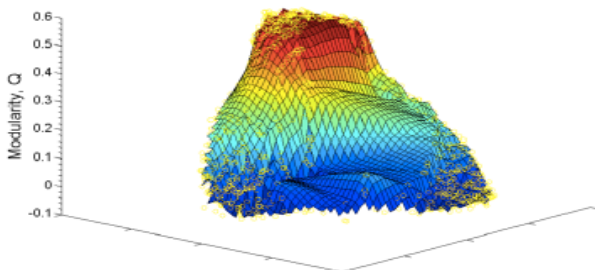
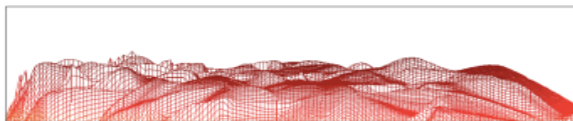
# MODULARITY RESOLUTION LIMITE



For  $m = 3$   $Max_Q = 0.675$  while natural partition has  $Q = 0.65$



# MODULARITY DISTRIBUTION



# PROPAGATION-BASED APPROACHES

---

## Algorithm 1 Label propagation

---

**Require:**  $G = \langle V, E \rangle$  a connected graph,

- 1: Initialize each node with unique label  $l_v$
- 2: **while** Labels are not stable **do**
- 3:   **for**  $v \in V$  **do**  
       $l_v = l_{|\Gamma^l(v)|}$  /\* random tie-breaking \*/
- 4:   **end for**
- 5: **end while**
- 6: **return** communities from labels

---

$\Gamma^l(v)$  : set of neighbors having label  $l$

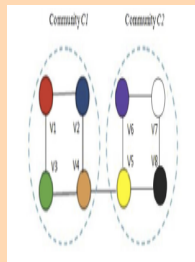
# LABEL PROPAGATION

## Advantages

- ▶ Complexity :  $\mathcal{O}(m)$
- ▶ Highly parallel

## Disadvantages

- ▶ No convergence guarantee, oscillation phenomena
- ▶ Low robustness  
*Different runs yields very different community structure due to randomness*



# LABEL PROPAGATION: CONVERGENCE ENHANCEMENT

- ▶ Asynchronous label update, [RAK07]
- ▶ Semi-synchronous label update (graph coloring + propagation by color) [CG12].
- ▶ Making stability even worse !
- ▶ Harden parallel implementations.

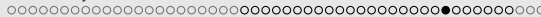
# ROBUST LABEL PROPAGATION

- ▶ *Label hop attenuation* [LHLC09]
- ▶ *Balanced label propagation* [SB11].
- ▶ **Multiplex approach !**  
*adding new neighborhood similarity based relationships between adjacent nodes → multiplex network*
- ▶ **Ensemble clustering** → Communities core [OGS10, SG12, LF12].

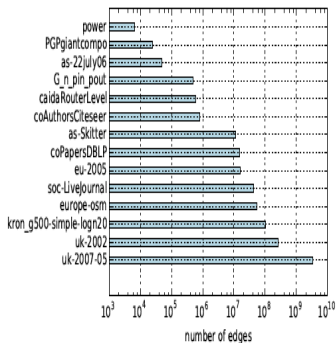


# LABEL PROPAGATION PREPROCESSING (EPP)

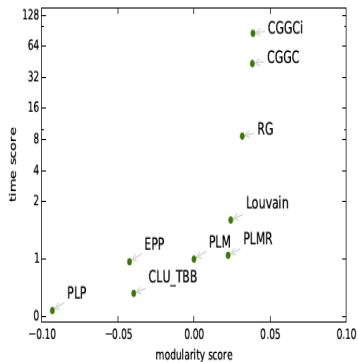
- 1 Apply an ensemble clustering on results of basic Parallel LP
- 2 Coarse the graph according to obtained communities
- 3 Apply a high quality community detection algorithm on coarsened graph.
- 4 Expand obtained results to the initial graph.



# EVALUATIONS



benchmark networks



Pareto front

# SEED-CENTRIC ALGORITHMS

(KANAWATI, SCSM'2014)

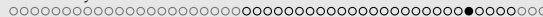
---

**Algorithm 2** General seed-centric community detection algorithm

---

**Require:**  $G = \langle V, E \rangle$  a connected graph,

- 1:  $\mathcal{C} \leftarrow \emptyset$
  - 2:  $S \leftarrow \text{compute\_seeds}(G)$
  - 3: **for**  $s \in S$  **do**
  - 4:    $C_s \leftarrow \text{compute\_local\_com}(s, G)$
  - 5:    $\mathcal{C} \leftarrow \mathcal{C} + C_s$
  - 6: **end for**
  - 7: **return**  $\text{compute\_community}(\mathcal{C})$
-



# SEED-CENTRIC APPROACHES

Table: Characteristics of major seed-centric algorithms

Algorithm	Seed Nature	Seed Number	Seed selection	Local Com.	Com. computation
Leaders-Followers [SZ10]	Single	Computed	informed	Agglomerative	-
Top-Leaders [KCZ10]	Single	Input	Random	Expansion	-
PapadopoulosKVS12	Subgraph	Computed	Informed	Expansion	-
WhangGD13	Single	Computed	informed	Expansion	-
BollobasR09	Subgraph	Computed	informed	expansion	-
Licod [Kan11]	Set	Computed	Informed	Agglomerative	-
Yasca [?]	Single	Computed	Informed	Expansion	Ensemble clustering

# EXAMPLE: LICOD

(KANAWATI, SOCILACOMP'2011)

## Licod : the idea !

- 1 Compute a set of seeds that are likely to be leaders in their communities  
*Heuristic : nodes having higher centralities than their neighbors*
- 2 Each node in the graph ranks seeds in function of its own preference
- 3 Each node modify its preference vector in function of neighbor's preferences
- 4 iterate max times or till convergence.



# THE YASCA ALGORITHM

(KANAWATI, COCOON'2014)

## The algorithm

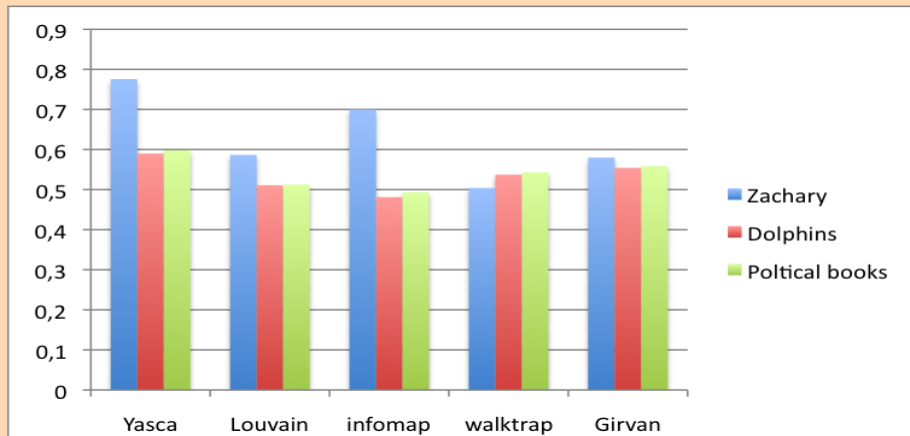
### Yet Another Seed-centric Community detection Algorithm

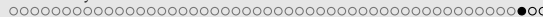
- 1 Compute a set of **diverse** seed nodes
- 2 For each node compute a bi-partition of the whole graphe based on local community identification.
- 3 Apply ensemble clustering to obtain a graph partition.

# YASCA: SOME RESULTS

## Comparative Results on benchmark networks : NMI

Select 15% of high central nodes and 15% of low central nodes





# EVALUATION METHODS

- ▶ **Topological criteria** : is it useful for applications ?
- ▶ Ground-truth comparaiison : hard to find on large-scale
- ▶ Task-oriented approaches : Clustering , Recommendation, link prediction, . . .

# TOPOLOGICAL EVALUATION CRITERIA

Quality of  $\mathcal{C} = \{S_1, \dots, S_i\}$  :

$$Q(\mathcal{C}) = \frac{\sum_i f(S_i)}{|\mathcal{C}|} \quad (2)$$

$f()$  is a single- community quality metric

4 types of quality functions

- Internal connectivity

- Externeal connectivity

- Hybrid functions

- Model-based functions: *the modularity*





# EXTERNAL CONNECTIVITY FUNCTIONS

$$\text{Expansion : } f(S) = \frac{C_S}{n_S}$$

$$\text{Cut : } f(S) = \frac{C_S}{n_S \times (N - n_S)}$$



# HYBRID FUNCTIONS

- ▶ **Conductance** :  $f(S) = \frac{C_S}{2m_S + C_S}$
- ▶ **MAX-ODF** : Out Degree Fraction :  $f(S) = \max_{u \in S} \frac{|\{(u,v) \in E, v \notin S\}|}{d(u)}$
- ▶ **AVG-ODF** :  $f(S) = \frac{1}{n_S} \times \sum_{u \in S} \frac{|\{(u,v) \in E, v \notin S\}|}{d(u)}$



## GROUND-TRUTH BASED EVALUATION

- ▶ Principle : Computing a *similarity* between obtained partition and a ground-truth one.
- ▶ Ground-truth :
  - Expert defined.
  - Model-generated
- ▶ Two types of metrics :
  - Agreement-based metrics: purity, rand, ARI
  - Information theory metrics: MI, NMI, . . . , etc.



# PURITY

- ▶  $\mathcal{P} = \{p_1, p_2, \dots, p_k\}$  : a computed partition
- ▶  $\mathcal{R} = \{r_1, r_2, \dots, r_n\}$  : ground-truth partition
- ▶  $\text{purity}(\mathcal{P}, \mathcal{R}) = \frac{1}{|V|} \sum_{j=1}^k \max_i (|p_k \cap r_i|)$



## RAND INDEX

- ▶  $\mathcal{P} = \{p_1, p_2, \dots, p_k\}$  : a computed partition
- ▶  $\mathcal{R} = \{r_1, r_2, \dots, r_n\}$  : ground-truth partition
- ▶  $a$  : # **pairs of nodes** in a same community according to  $\mathcal{P}$  and  $\mathcal{R}$
- ▶  $b$  : # **pairs of nodes** in a same community according to  $\mathcal{P}$  and in different communities in  $\mathcal{R}$
- ▶  $c$  : # **pairs of nodes** in a different communities according to  $\mathcal{P}$  and in same community in  $\mathcal{R}$
- ▶  $d$  : # **pairs of nodes** in a different communities in  $\mathcal{P}$  and  $\mathcal{R}$
- ▶
- ▶  $rand(\mathcal{P}, \mathcal{R}) = \frac{a+d}{a+b+c+d}$
- ▶  $ARI = \frac{C_n^2(a+d) - [(a+b)(a+c) + (c+d)(b+d)]}{(C_n^2)^2 - [(a+b)(a+c) + (c+d)(b+d)]}$
- ▶  $E(ARI) = 0$
- ▶ **rappel** :  $C_n^k = \frac{n!}{k!(n-k)!}$











# COMMUNITY DETECTION IN MULTIPLEX NETWORKS

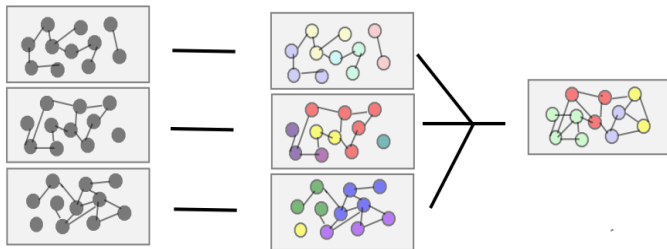
## Approaches

- 1 Transformation into a monoplex community detection problem**
  - ▶ Layer aggregation approaches.
  - ▶ Ensemble clustering approaches
  - ▶ Hypergraph transformation based approaches
- 2 Generalization of monoplex oriented algorithms to multiplex networks.**
  - ▶ Generalized-modularity optimization
  - ▶ **Seed-centric approaches**

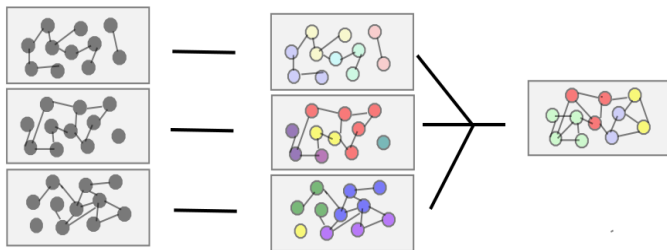




# ENSEMBLE CLUSTERING APPROACHES



# ENSEMBLE CLUSTERING APPROACHES



## Ensemble Clustering

Strehl2003

- ▶ CSPA: Cluster-based Similarity Partitioning Algorithm
- ▶ HGPA: HyperGraph-Partitioning Algorithm
- ▶ MCLA: Meta-Clustering Algorithm
- ▶ ...

# K-UNIFORM HYPERGRAPH TRANSFORMATION

## KIVELA2013MULTILAYER

### Principle

- ▶ A  $k$ -uniform hypergraph is a hypergraph in which the cardinality of each hyperedge is exactly  $k$
- ▶ Mapping a multiplex to a **3-uniform hypergraph**  $\mathcal{H} = (\mathcal{V}, \mathcal{E})$  such that :
  - $\mathcal{V} = V \cup \{1, \dots, \alpha\}$
  - $(u, v, i) \in \mathcal{E}$  if  $\exists l : A_{uv}^{[l]} \neq 0, u, v \in V, i \in \{1, \dots, \alpha\}$
- ▶ Apply community detection approaches in Hypergraphs (Ex. tensor factorization approaches)

















# RESULTS: MUXLICOD

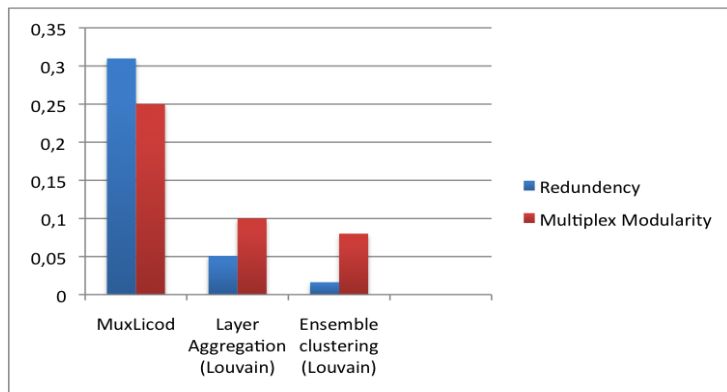


Figure: DBLP 3-layer multiplex











