



BIG DATA ANALYTICS

A Social Network Approach



Andry Alamsyah
Social Computing and Big Data Research Group
Fakultas Ekonomi dan Bisnis

Andry Alamsyah

Fakultas Ekonomi dan Bisnis

Universitas Telkom

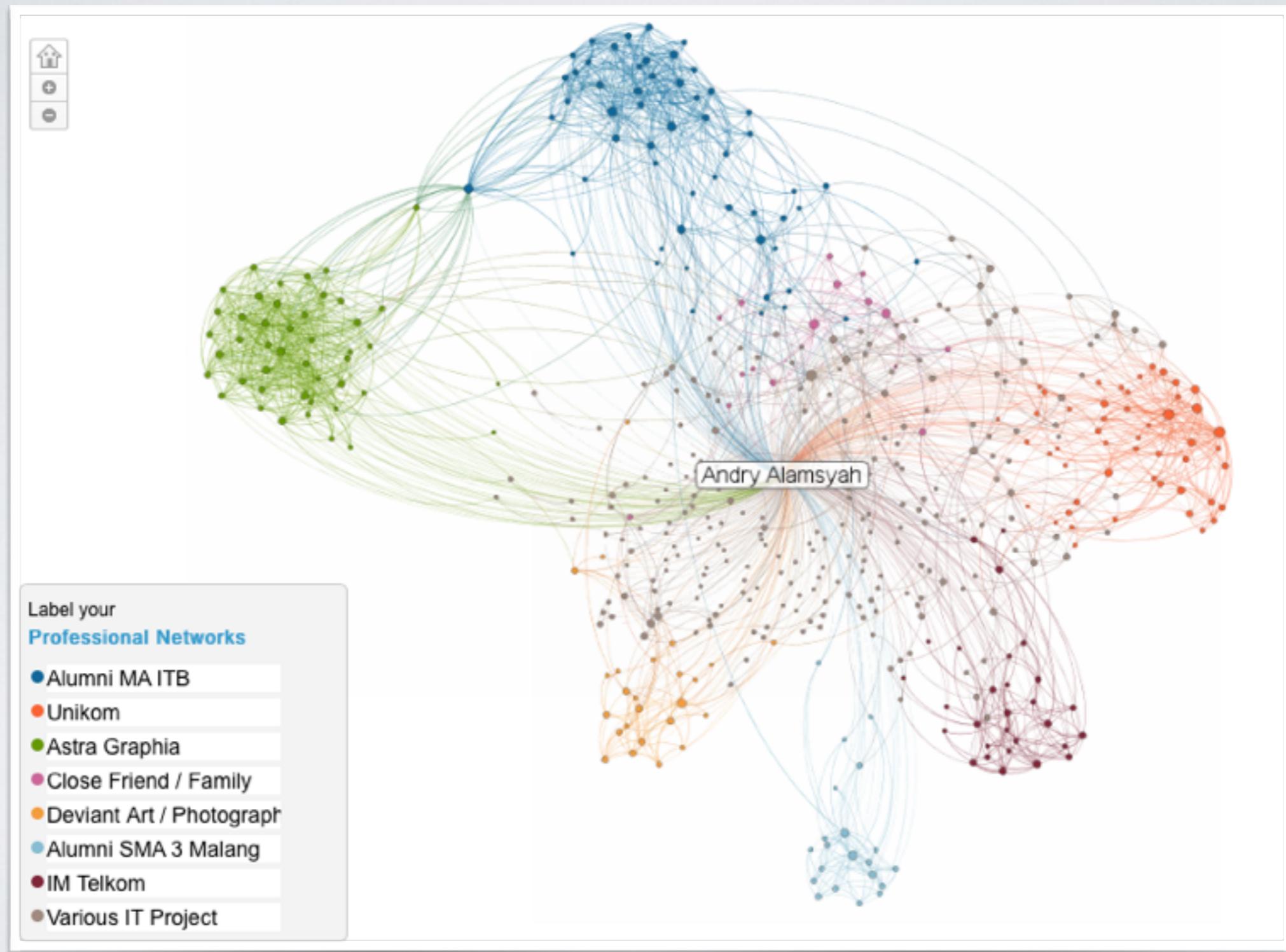
Research Field :

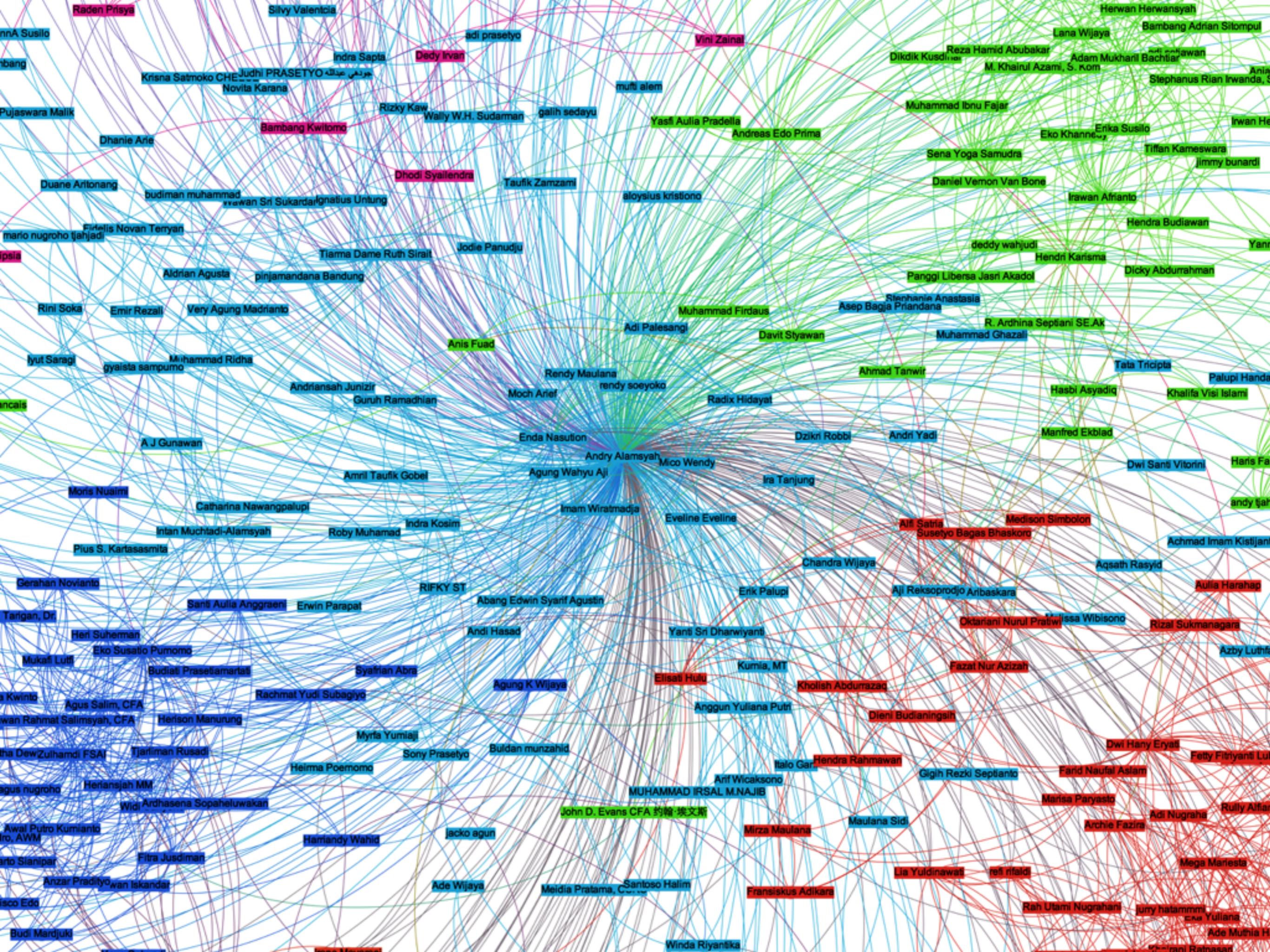
Social Network, Complex Network / Network Science, Social Computing, Data Analytics, Data Mining, Big Data, Graph Theory, Content Business, Data Business, ICT Business



andry.alamsyah@gmail.com
andrya.staff.telkomuniversity.ac.id
telkomuniversity.academia.edu/andryalamsyah
researchgate.net/profile/Andry_Alamsyah
linkedin.com/in/andry.alamsyah
twitter.com/andrybrew

WHO AM I ?





LARGE SCALE DATA



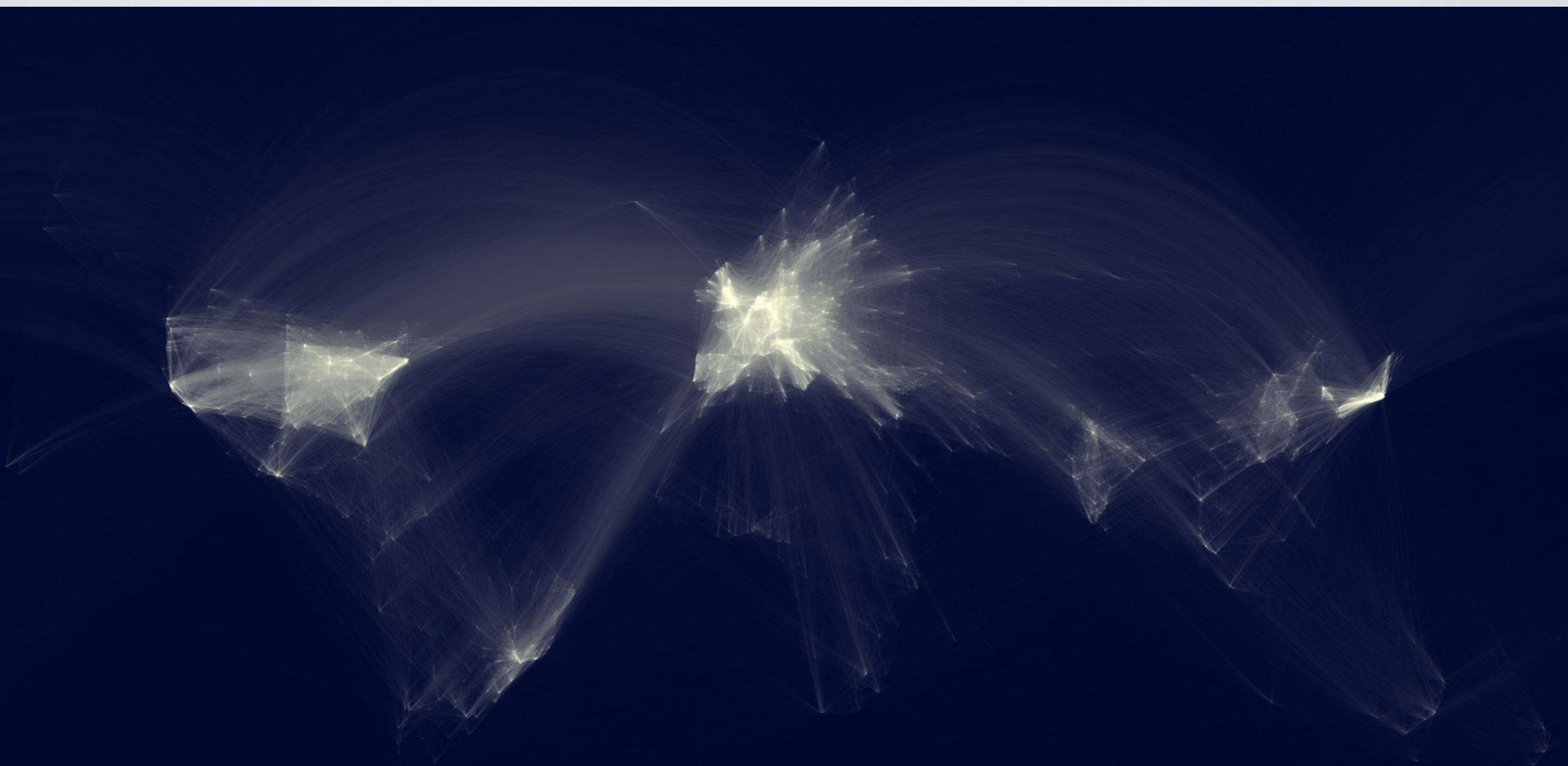
facebook

December 2010

LARGE SCALE DATA



LARGE SCALE DATA



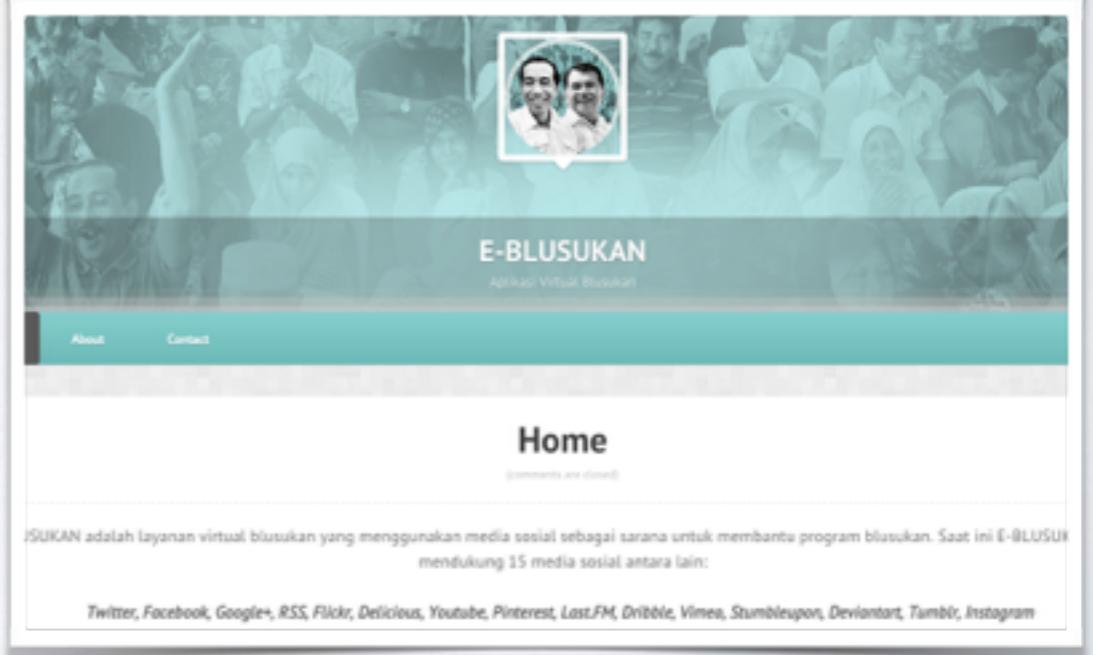
Map of scientific collaborations from 2005 to 2009

Computed by Olivier H. Beauchesne @ Science-Metrix, Inc.

Data from Scopus, using books, trade journals and peer-reviewed journals

STORY / PHENOMENON

- BIG DATA leads to Social Computing
(*Quantification of Individual / Social Behaviour*)
- Social Network Data / Conversation are widely available
- Social Network voices represent public voice become '**Big**' concern (references)
- The Need of Real-Time Analytic (OLAP)
- The Need of Powerful Metric for Social Network / Big Data



STORY / PHENOMENON

- There are many aspect of Big Data research, but too little resource, too little talent
- Same business objective, but increase effectiveness on top of current services
- Problem with Legacy Methodology approach using Questionnaire/Interviews/Surveys (ok with small scale data , expensive and took longer time for large scale data, accuracy issues)



INDUSTRY EFFORTS

Method



METHODS COMPARISON IN SOCIAL SCIENCE

LEGACY

DATA ANALYTICS

Confirmative

Explorative (Predictive)

Small Data Set

Larga Data Set

Small Number of Variable

Large Number of Variable

Deductive (no predictions)

Inductive

Numeric Data

Numeric and Non-Numeric Data

Clean Data

Data Cleaning

BIG DATA STATE OF THE ART

Computation Related

Processing / Computation	Storage	Analytics Tools
<ul style="list-style-type: none">• Hadoop• Nvidia CUDA• Twitter Storm• Bulk Synchronous Parallel Processing• GraphLab• Disk-Based Graph Processing	<ul style="list-style-type: none">• neo4j• Titan• HDFS	<ul style="list-style-type: none">• MLPACK• Mahout

Methodology / Analytics Related

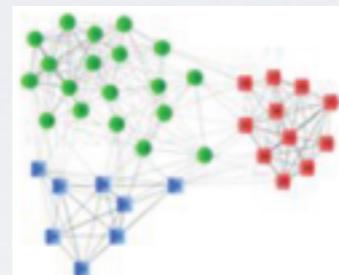
modelling, descriptions, predictions, optimisation and simulation

BIG DATA ANALYTICS CONSTRUCTOR

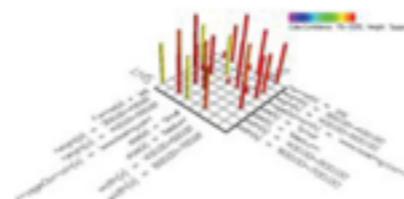
Social Network

networks
tie-strength
key players
cohesion

Data Mining



Clustering



Association Rules



Classification & Regression

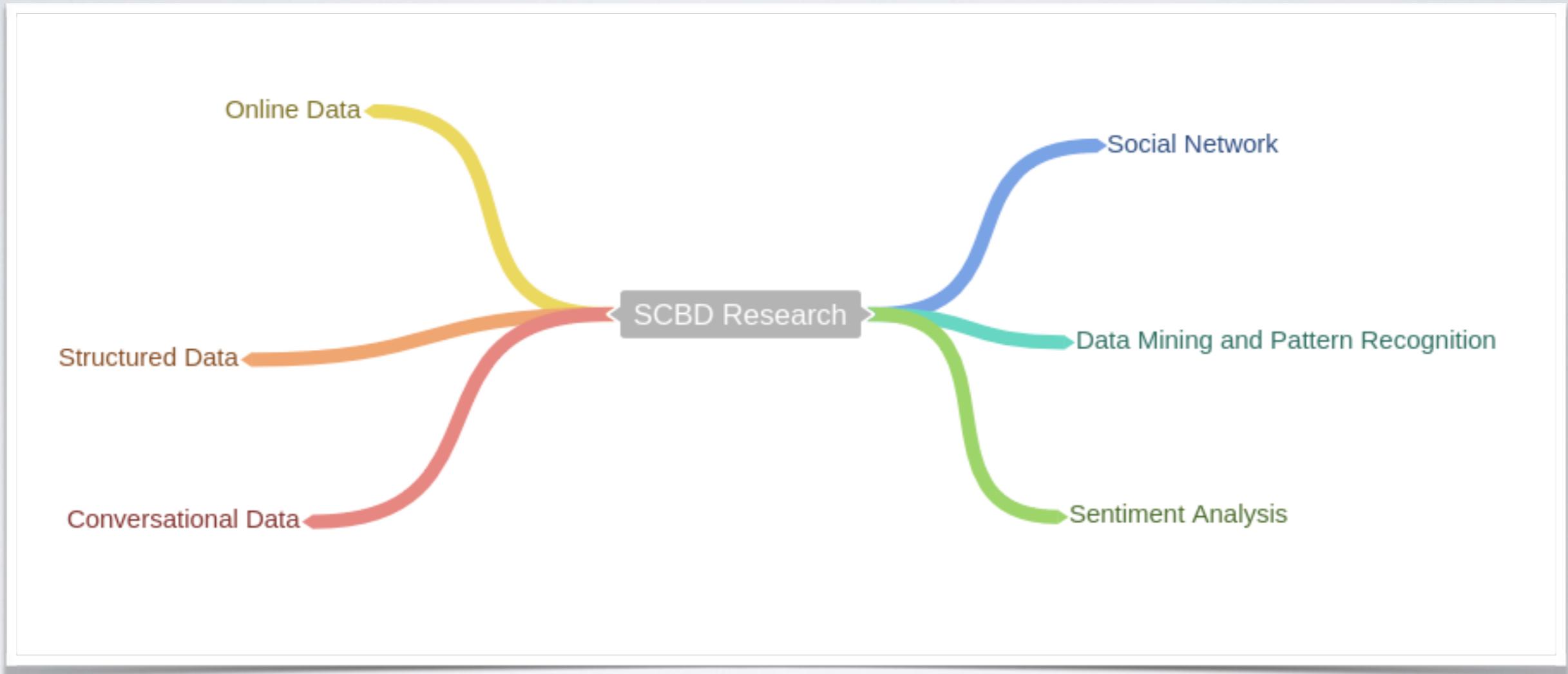


Anomaly Detection

Sentiment Analysis

keyword spotting
lexical affinity
statistical methods
concept-level technique

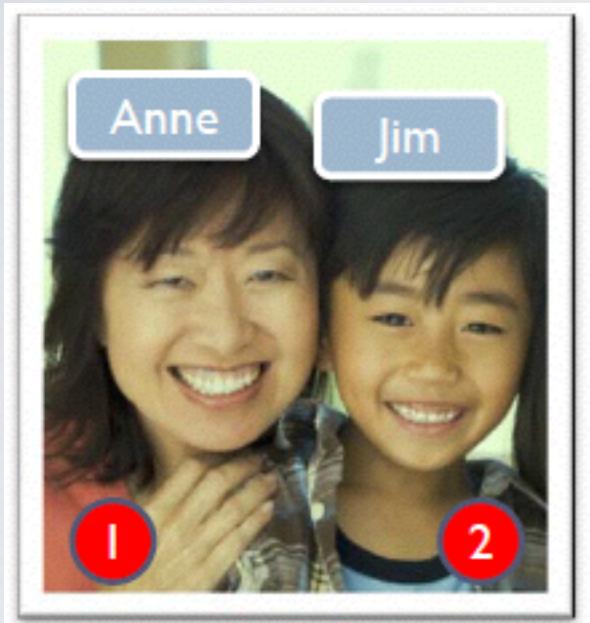
RESEARCH ROADMAP



GOAL : descriptions, predictions, optimisation and simulation

area : marketing, communications, knowledge
management, operations, finance, etc

SOCIAL NETWORK MODEL



Can we study their interactions as a network ?

Communication

Anne: Jim, tell the Murrays they're invited

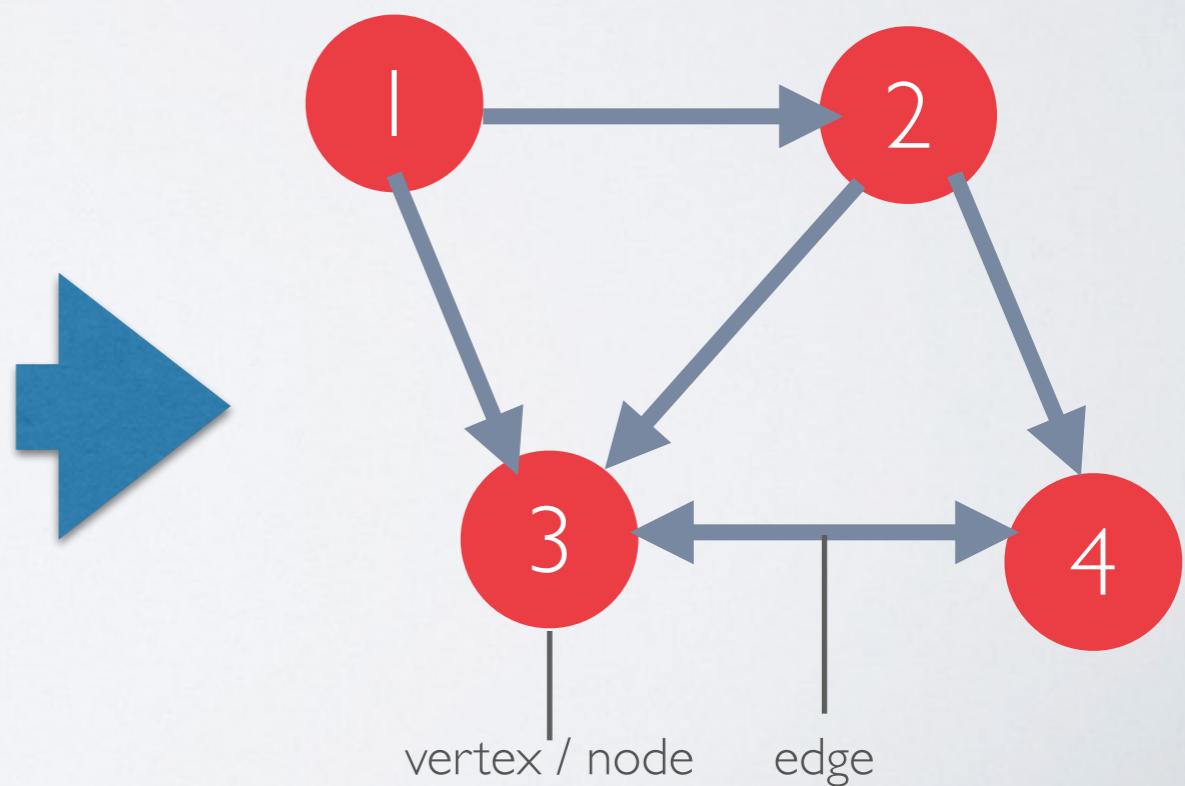
Jim: Mary, you and your dad should come for dinner!

Jim: Mr. Murray, you should both come for dinner

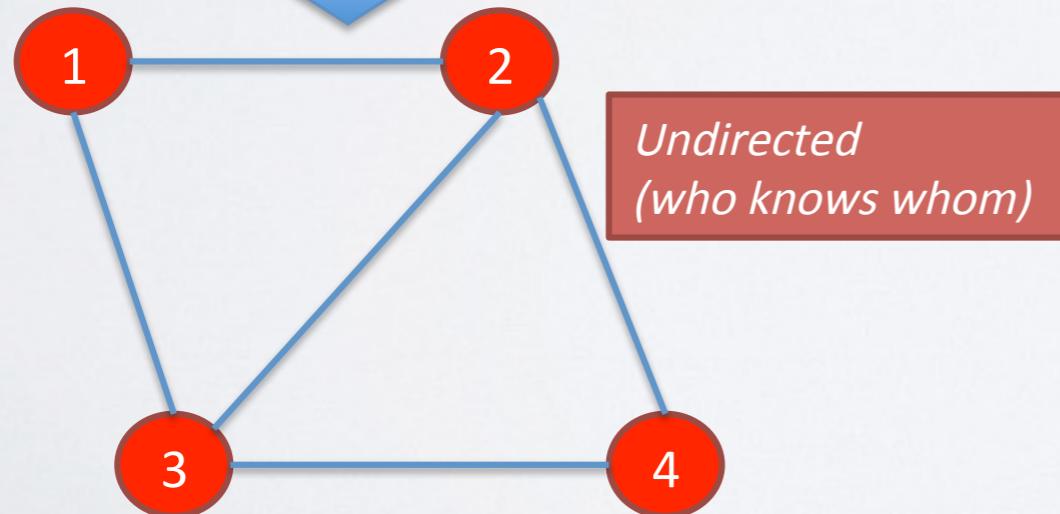
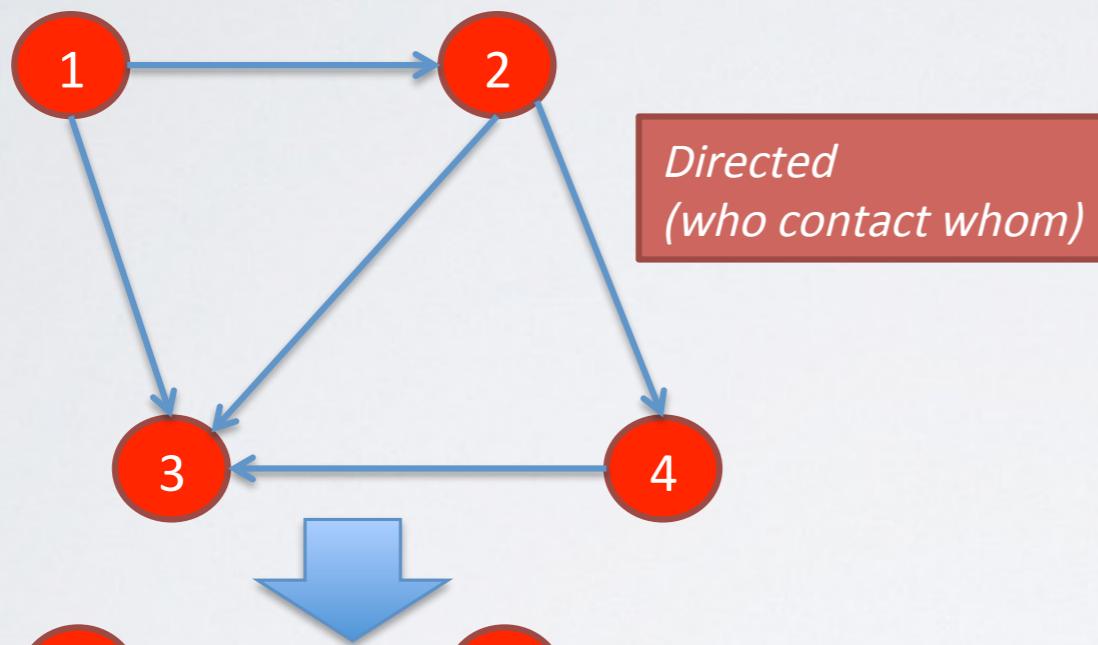
Anne: Mary, did Jim tell you about the dinner? You must come.

Mary: Dad, we are invited for dinner tonight

John: (to Anne) Ok, we're going, it's settled!



SOCIAL NETWORK MODEL



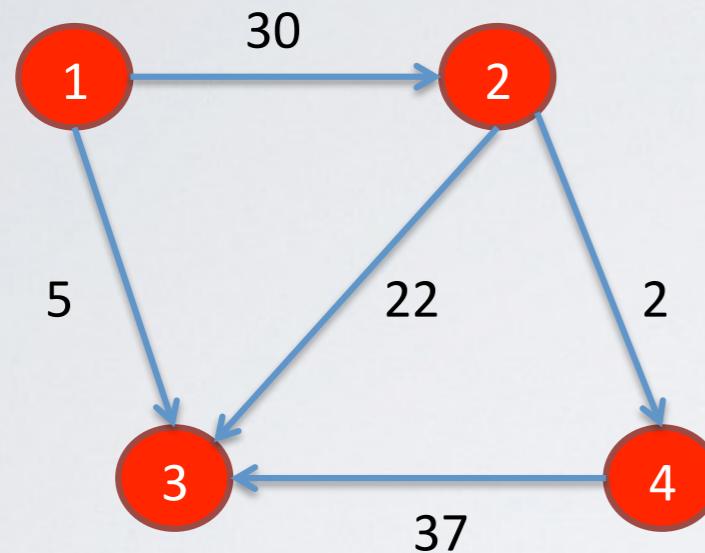
Edges List

Vertex	Vertex
1	2
1	3
2	3
2	4
4	3

Adjacency Matrix become symmetric

Vertex	1	2	3	4
1	-	1	1	0
2	1	-	1	1
3	1	1	-	0
4	0	1	0	-

TIE STRENGTH



Weight could be

- Frequency of interactions in period of observation
- Number of items exchanged in period
- Individual perceptions of strength of relationship
- Cost of communications or exchange, e.g. distance

Edges List

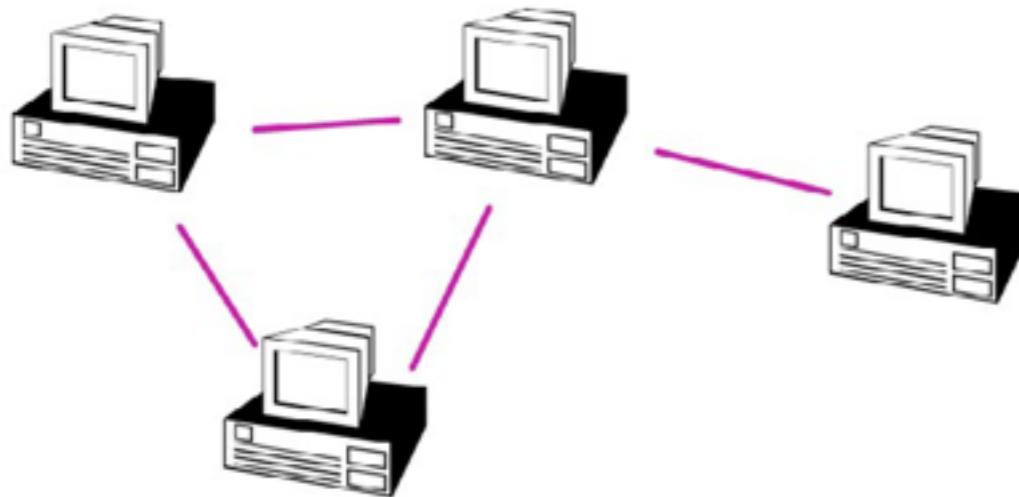
Vertex	Vertex	Weight
1	2	30
1	3	5
2	3	22
2	4	2
4	3	27

Adjacency Matrix (weight)

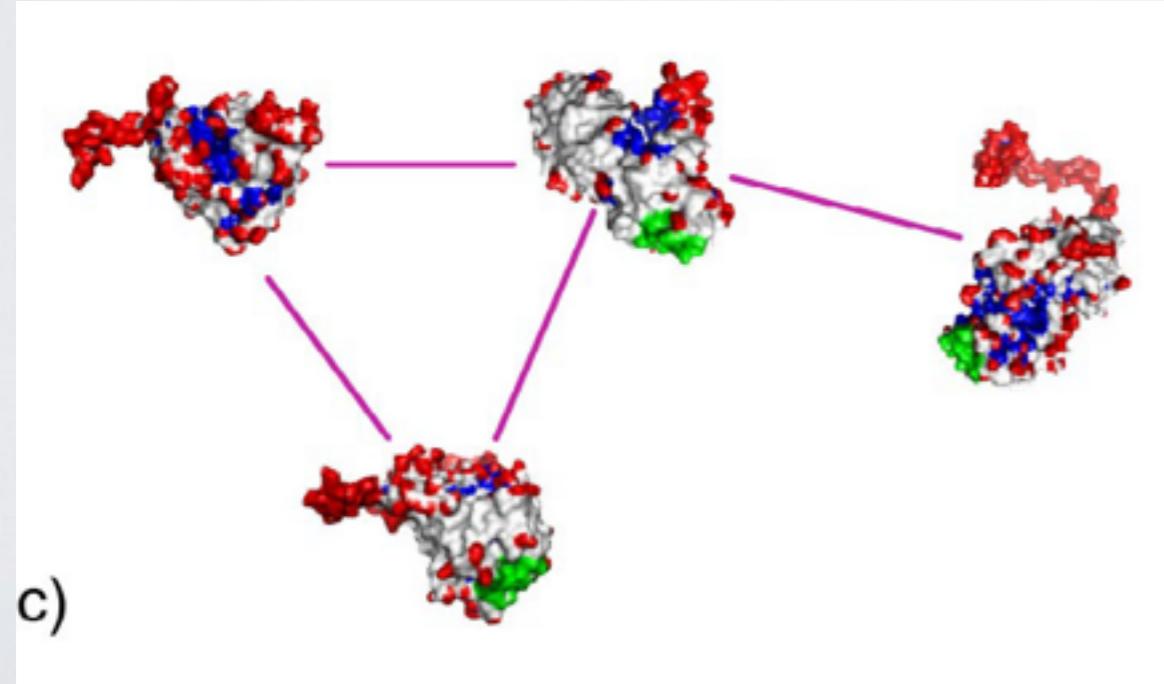
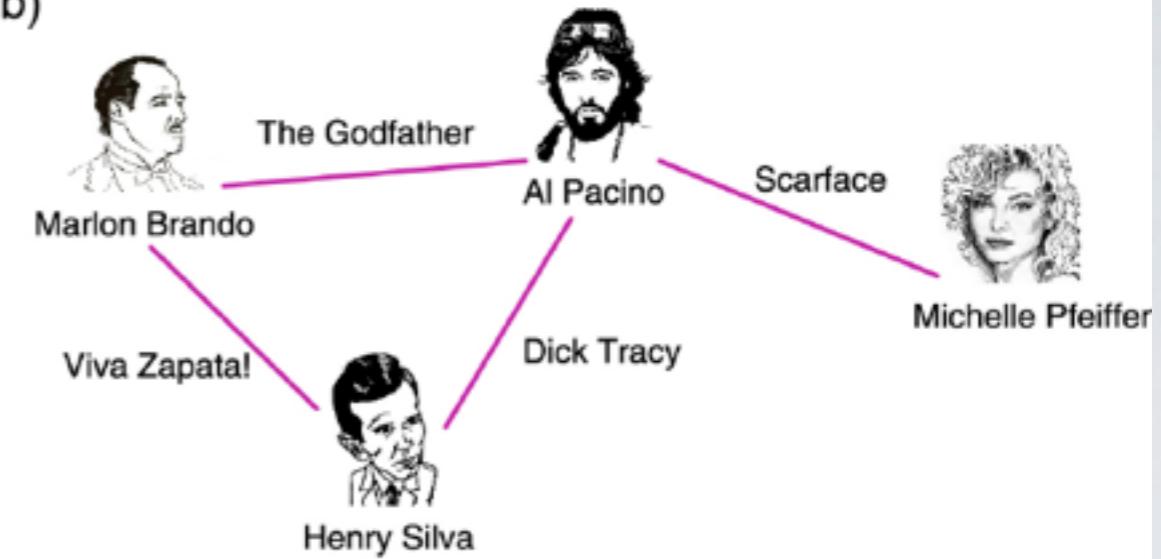
Vertex	1	2	3	4
1	-	30	5	0
2	30	-	22	2
3	5	22	-	37
4	0	2	37	-

NETWORK MODEL EXAMPLE

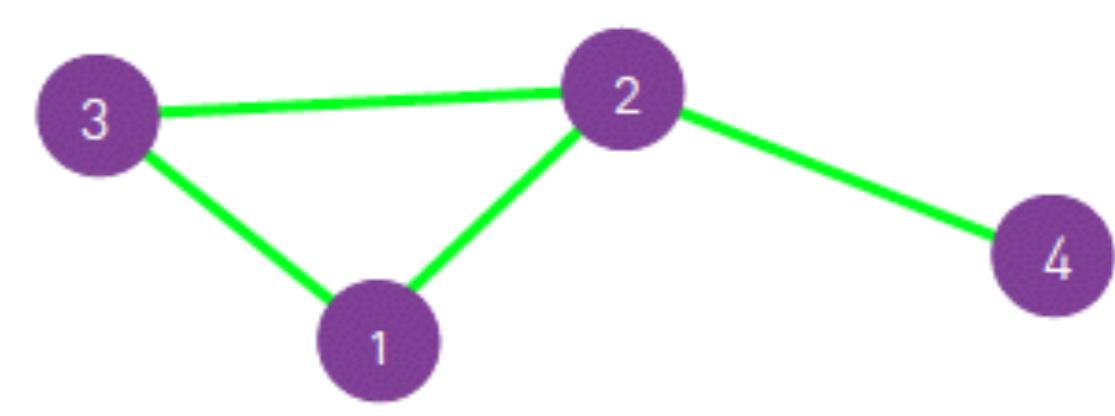
a)



b)

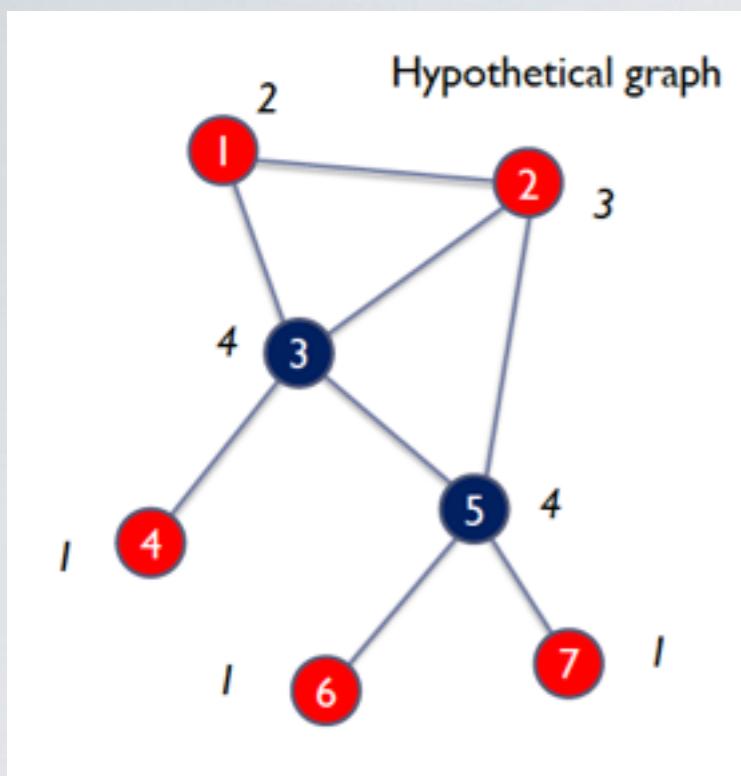


c)

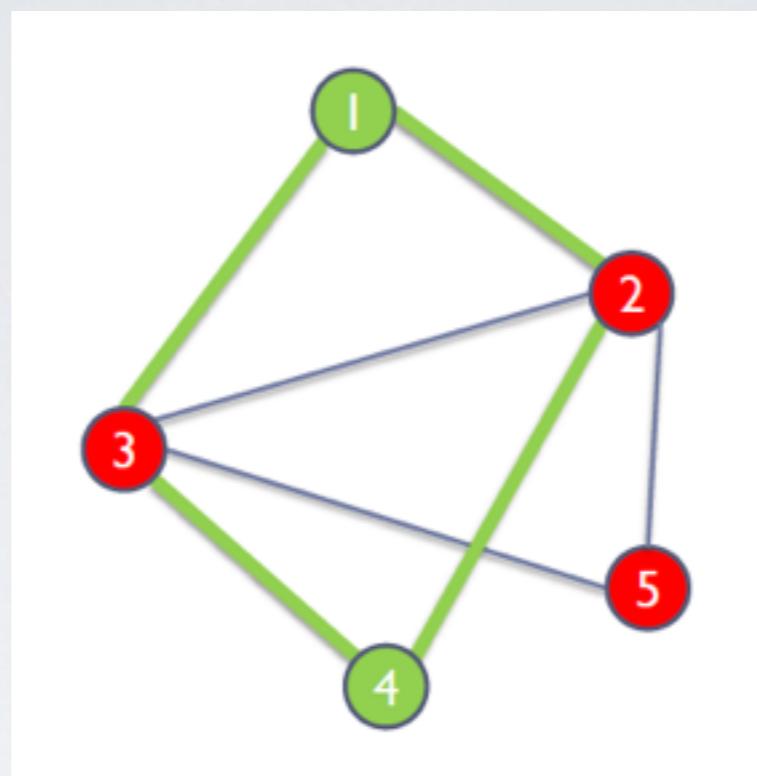


Different Network, Same Graph

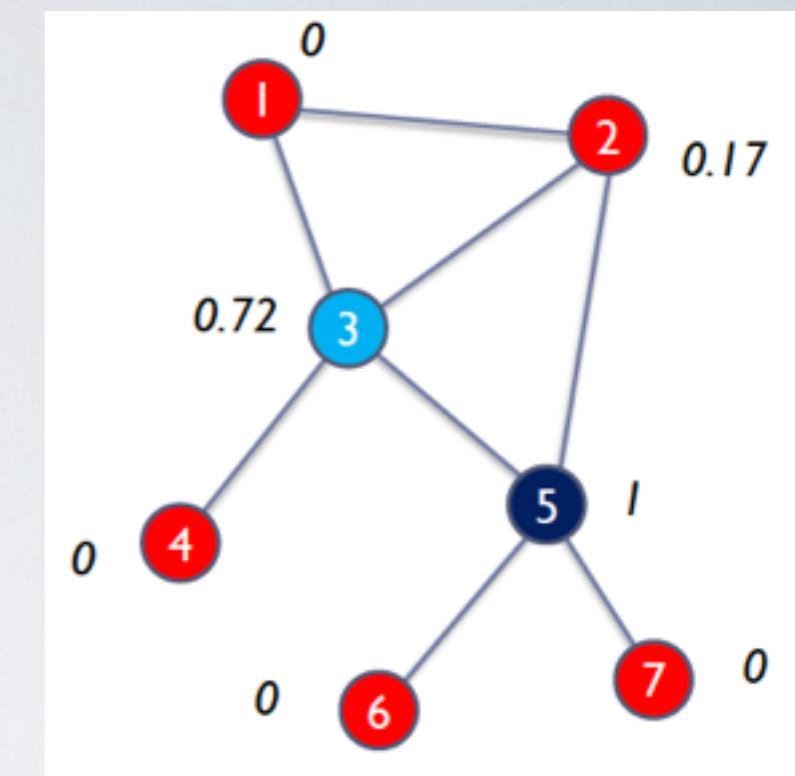
METRIK CENTRALITY



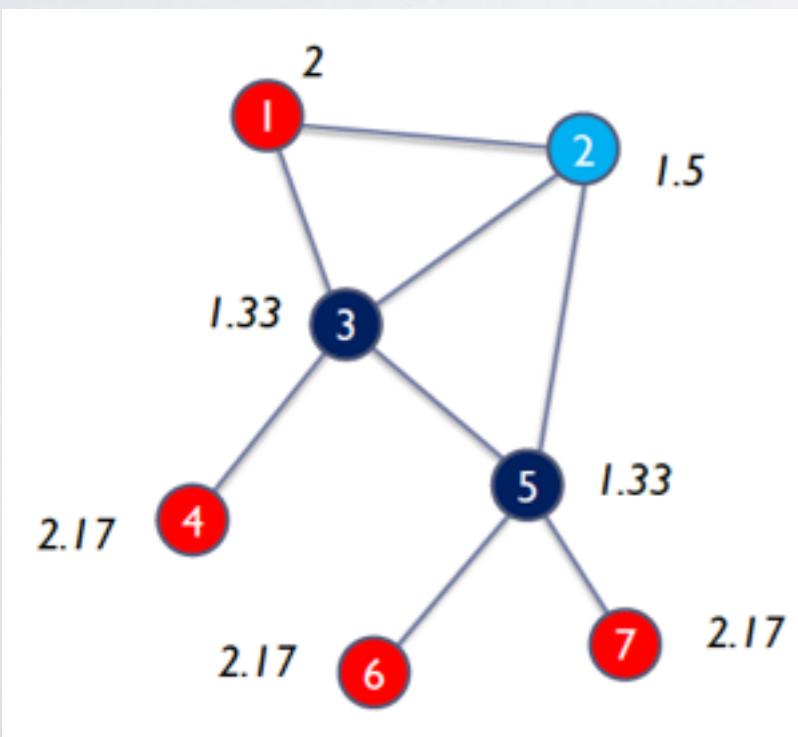
degree centrality



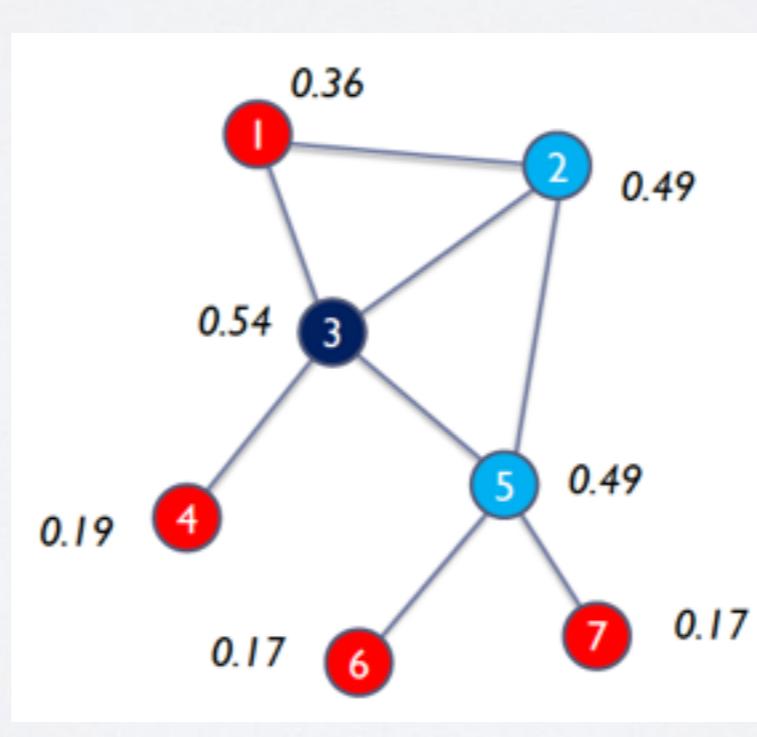
shortest path



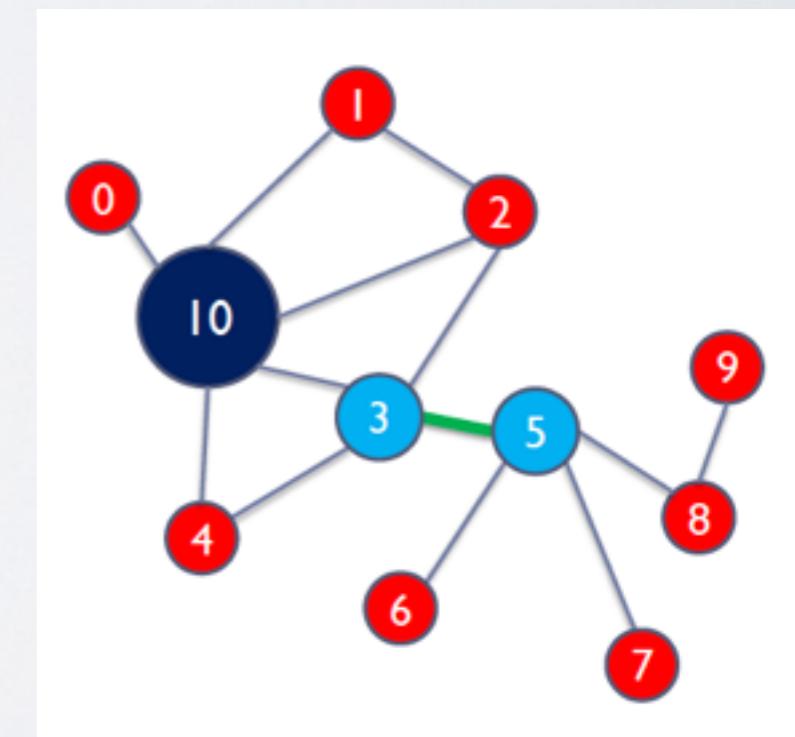
betweenness centrality



closeness centrality

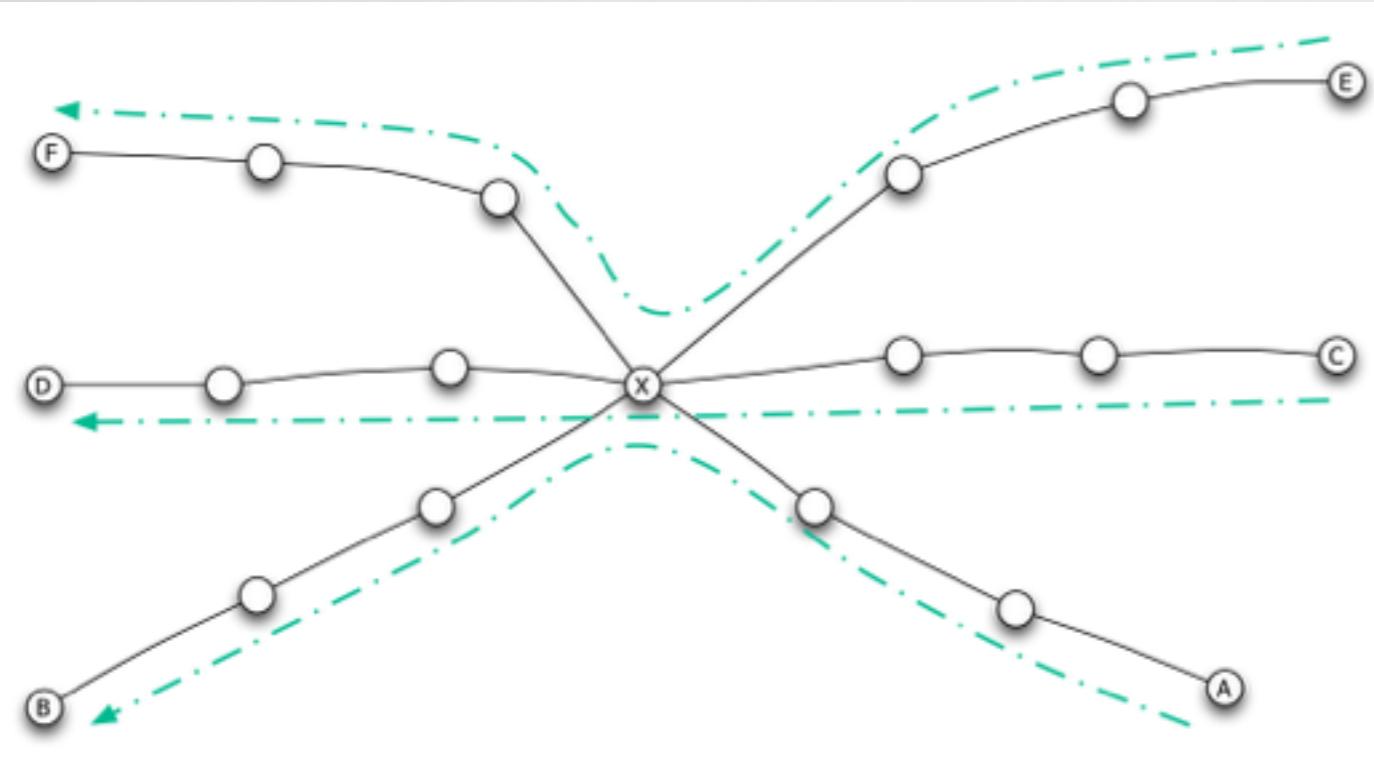


eigenvector centrality



set of key players

METRIK CENTRALITY

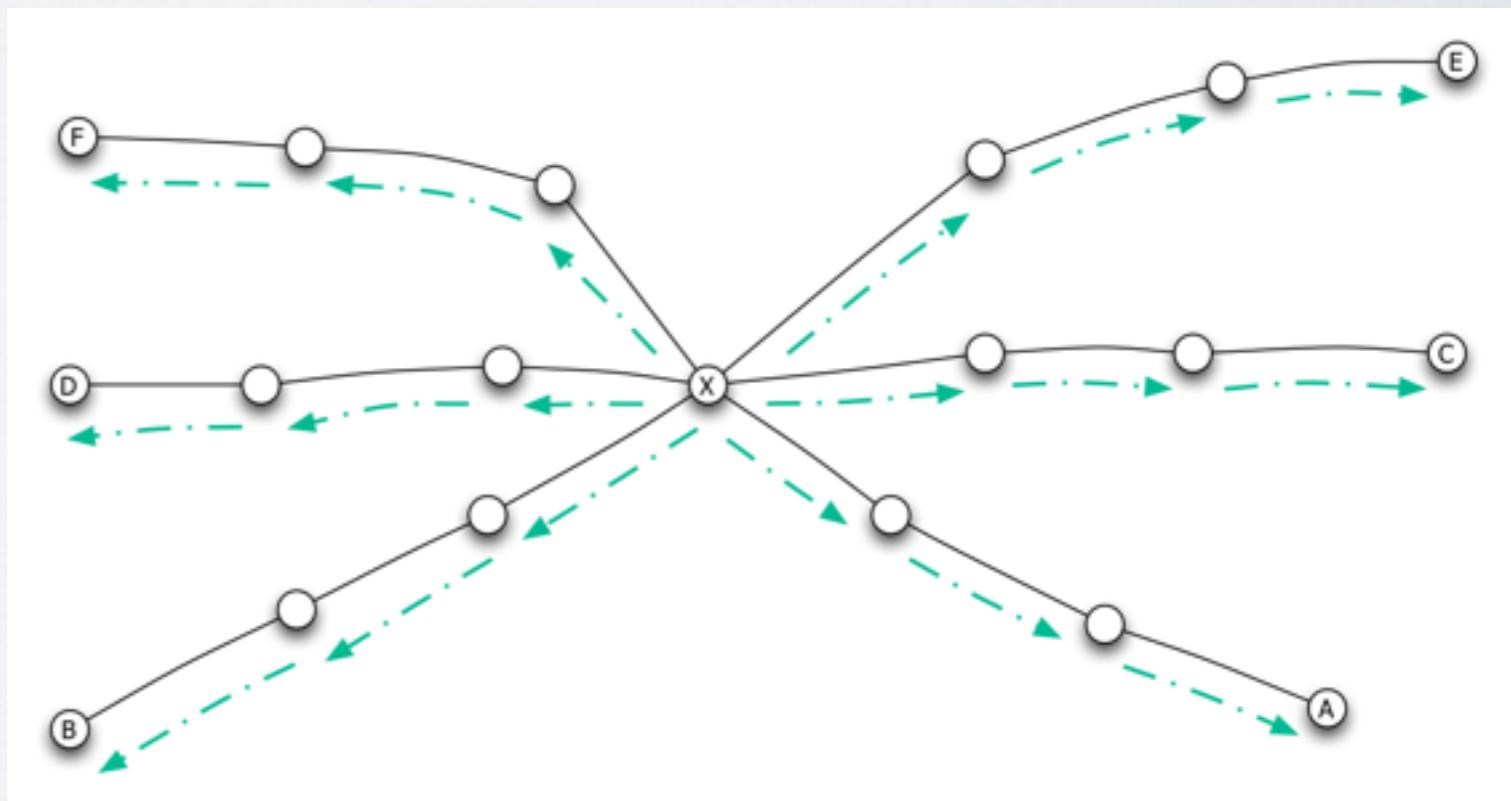


betweenness centrality

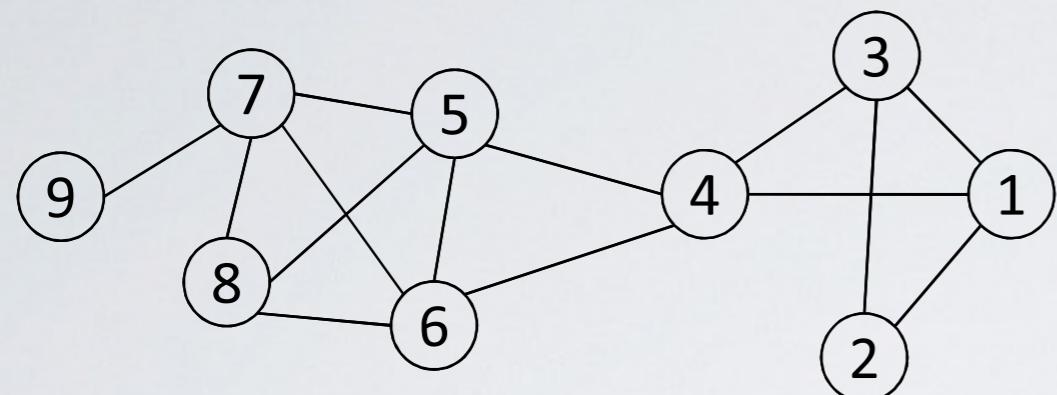
banyaknya **jalur terpendek** antar pasangan semua titik di jaringan, yang melewati satu titik yang diukur

closeness centrality

jarak titik yang diukur terhadap semua titik yang ada dalam jaringan



FORMULASI CENTRALITY



*kompleksitas waktu perhitungan metrik mencapai $O(n^3)$ dan kompleksitas ruang sebesar $O(n^2)$

Betweenness Centrality

$$C_B(i) = \sum_{s \neq i \neq t} \frac{\sigma_{st}(i)}{\sigma_{st}}$$

$$\sigma_{st}(4) / \sigma_{st}$$

	$s=1$	$s=2$	$s=3$
$t=5$	1/1	2/2	1/1
$t=6$	1/1	2/2	1/1
$t=7$	2/2	4/4	2/2
$t=8$	2/2	4/4	2/2
$t=9$	2/2	4/4	2/2

$$C_B(4) = 15$$

Closeness Centrality

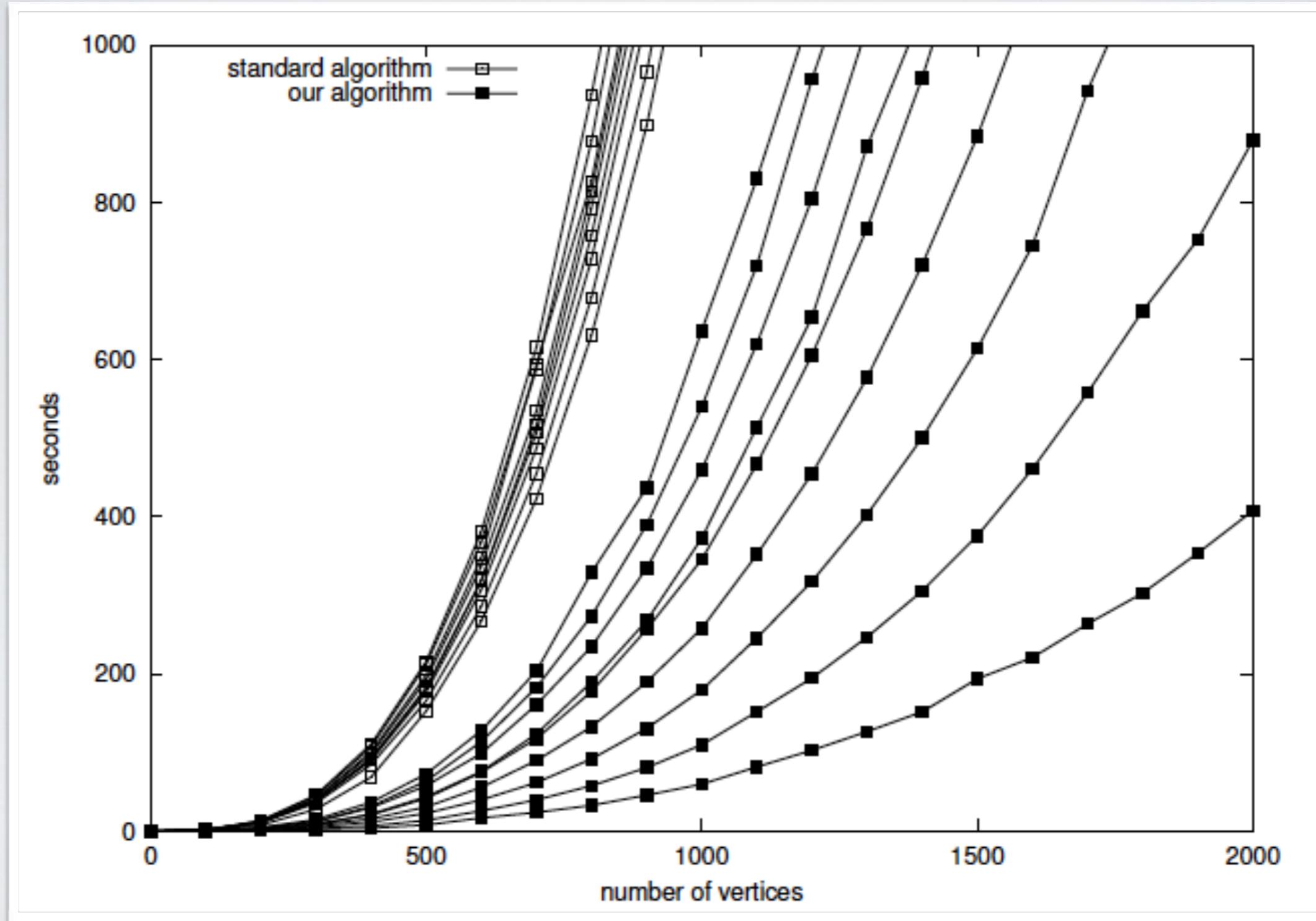
$$C_c(i) = \frac{n-1}{\sum_{j(\neq i)} d_{ij}}$$

	1	2	3	4	5	6	7	8	9
1	0	1	1	1	2	2	3	3	4
2	1	0	1	2	3	3	4	4	5
3	1	1	0	1	2	2	3	3	4
4	1	2	1	0	1	1	2	2	3
5	2	3	2	1	0	1	1	1	2
6	2	3	2	1	1	0	1	1	2
7	3	4	3	2	1	1	0	1	1
8	3	4	3	2	1	1	1	0	2
9	4	5	4	3	2	2	1	2	0

$$C_c(4) = 0,62$$

$$C_c(3) = 0,47$$

COMPUTATION PROBLEM



CENTRALITY INTERPRETATIONS

Centrality Measures

Degree

Betweenness

Closeness

Eigenvector

Interpretation I

how many people can this person reach directly ?

how likely is this person to be the most direct route between two people in the network ?

how fast can this person reach everyone in the network ?

how well is this person connected to other well-connected people ?

Interpretation 2

in network of investors and start up: how many start up has this company invested to ?

in network of knowledge: who is the employee through whom most of the confidential information is likely to flow ?

in network of information dissemination: how fast this word-of-mouth issue spread from this person to the rest of the network ?

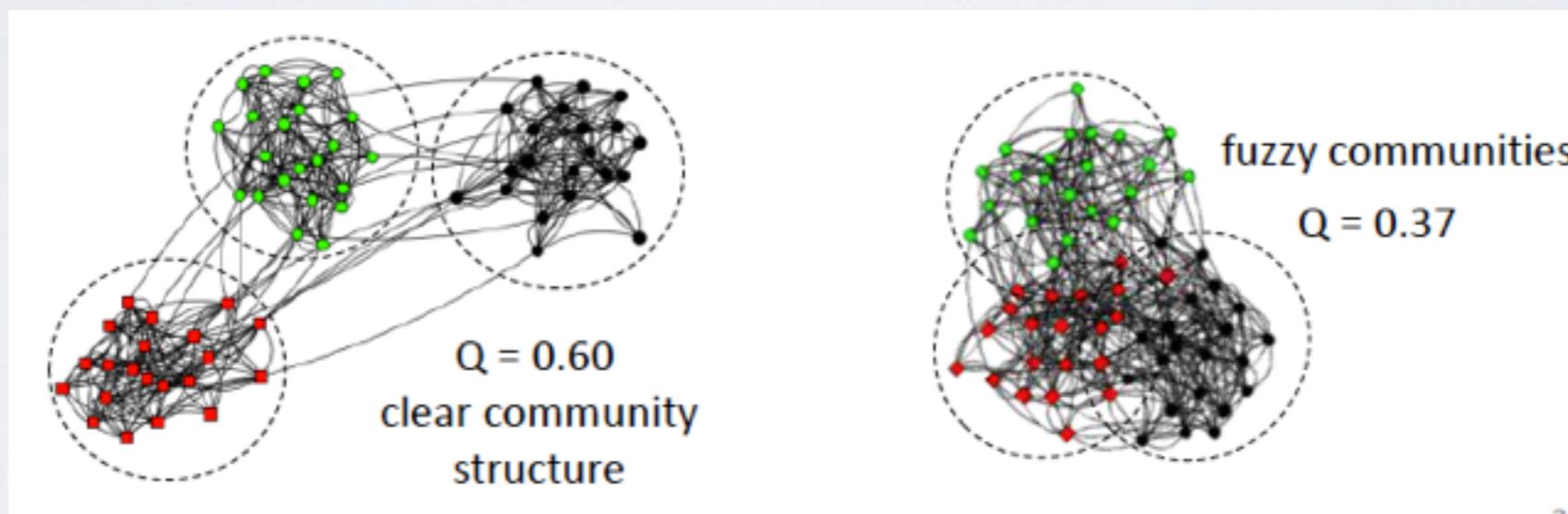
in network of paper citations: who is the author that is most cited by other well-cited authors ?

METRIK MODULARITY

$$Q = \frac{1}{2m} \sum_{ij} (A_{ij} - \frac{k_i k_j}{2m}) \delta(C_i, C_j)$$

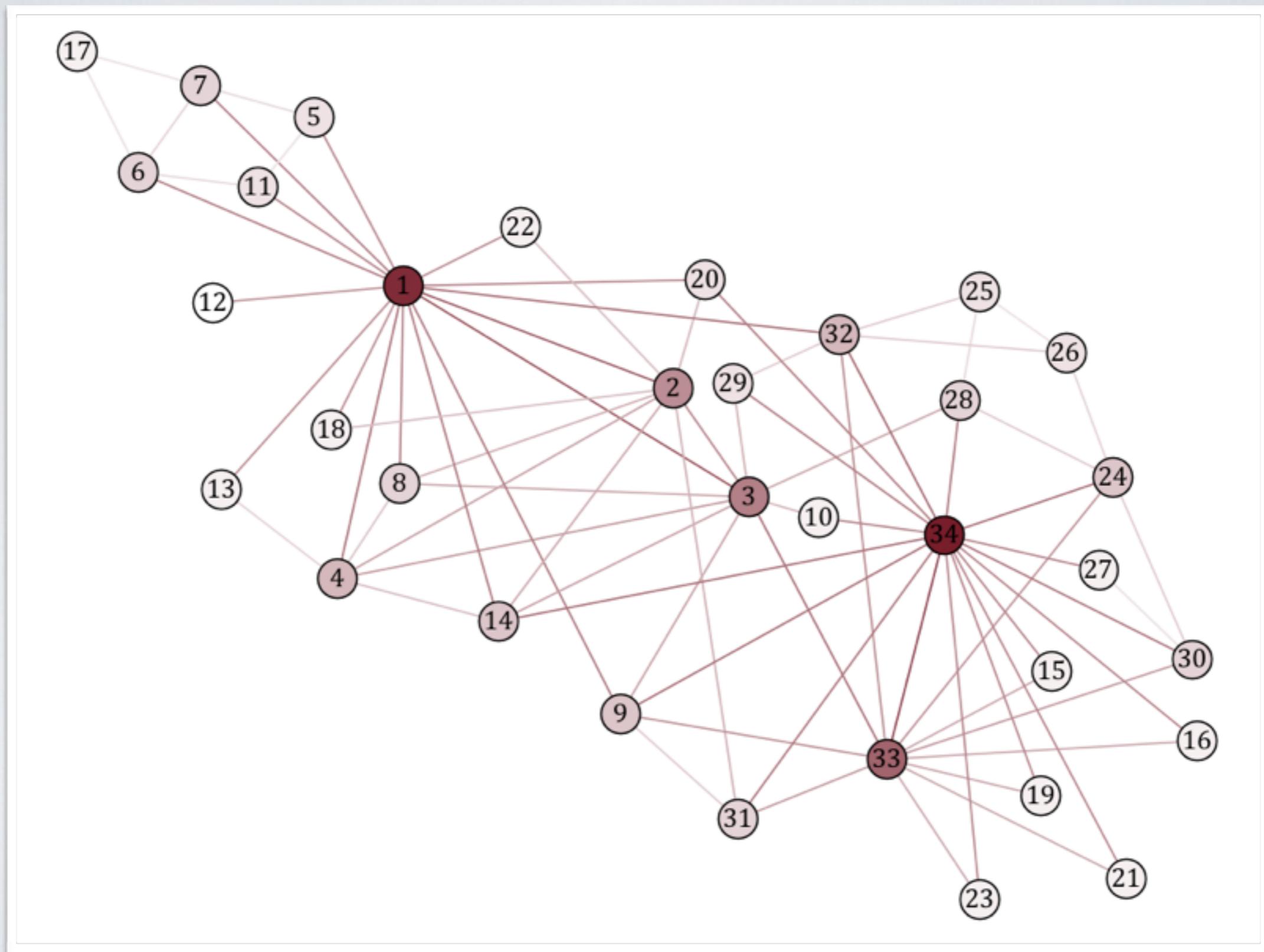
edges inside the community

expected number of edges if i, j places at random

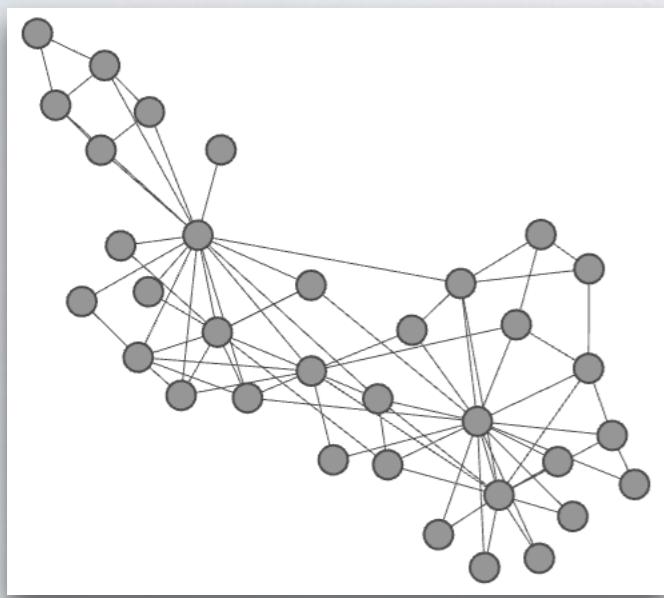


Metrik Modularity Kompleksitas Metrik
indeks kualitas partisi jaringan kompleksitas waktu perhitungan
menjadi komunitas $O(n^3)$

NETWORK MODEL EXAMPLE



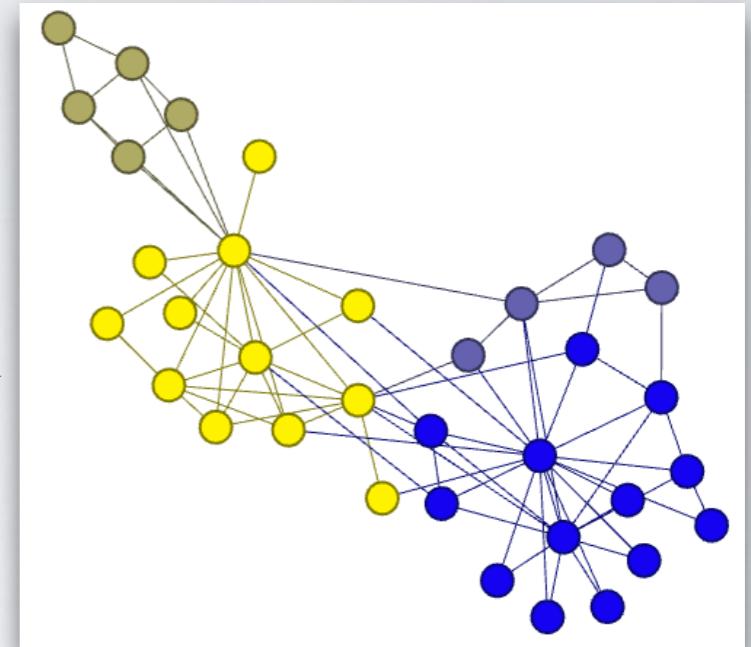
METRIC SIMULATION



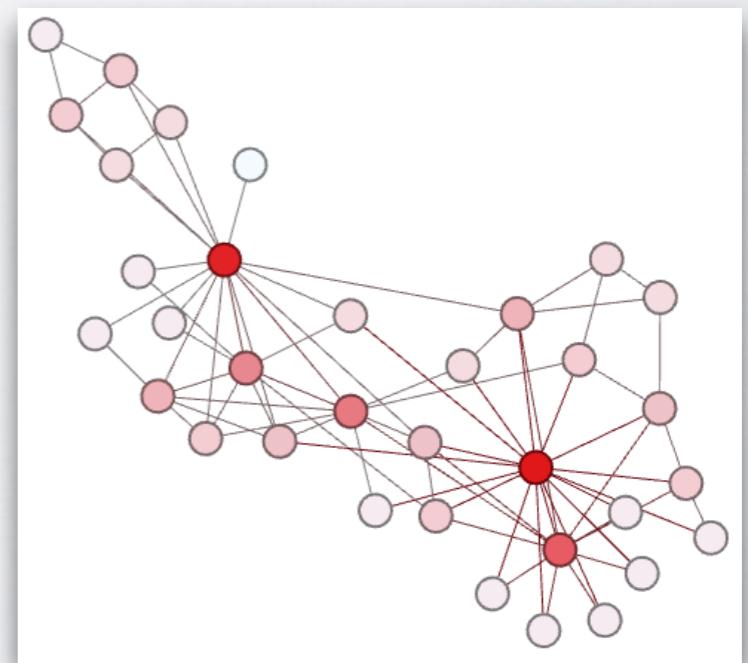
model network

modularity
community detection

centrality



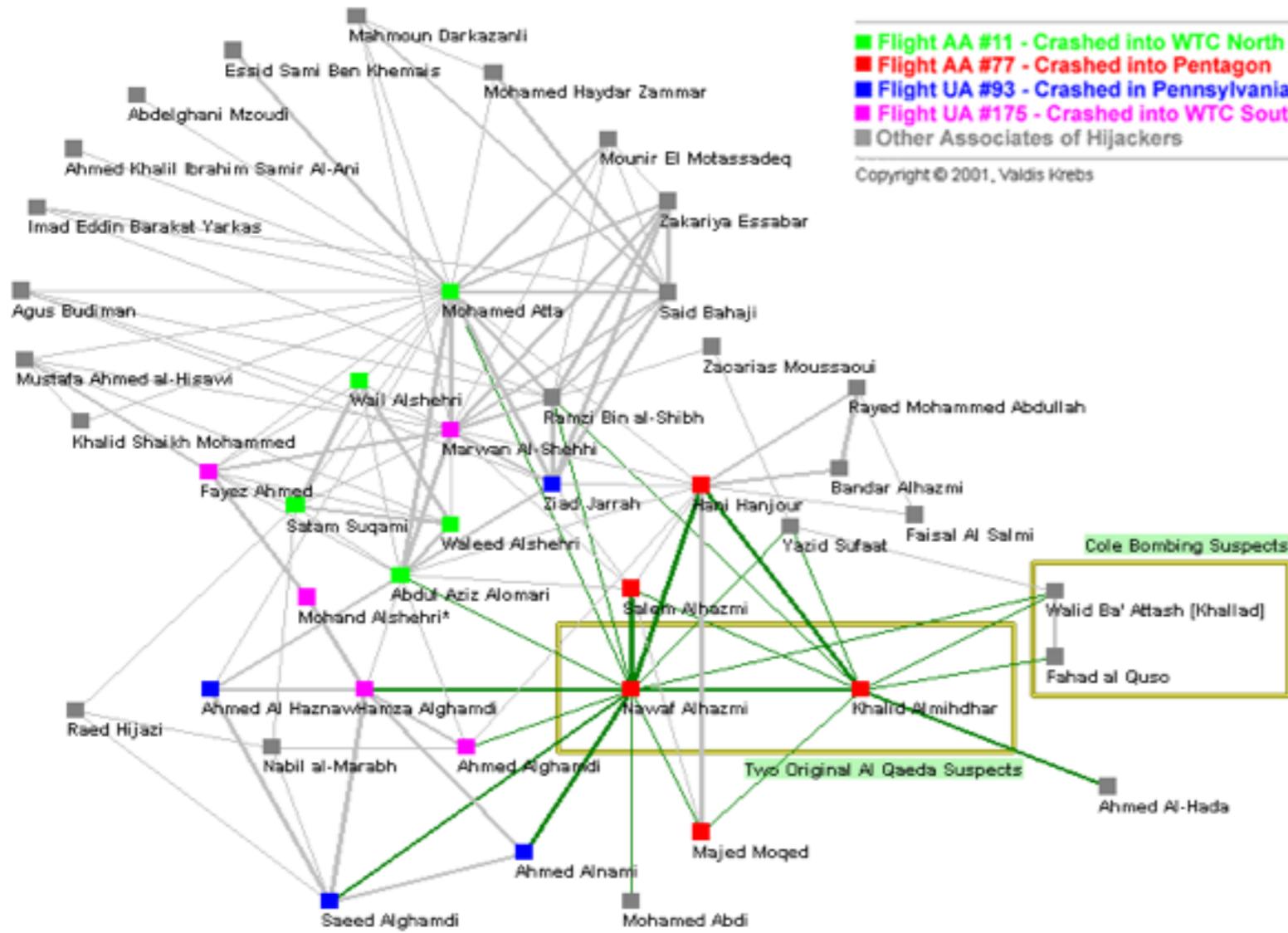
community detection result



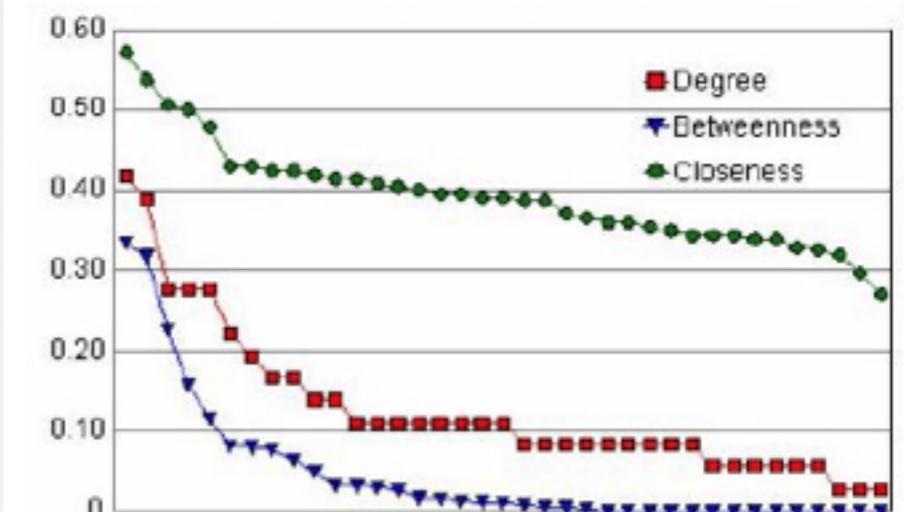
*karate club dataset
34 nodes dan 78 edges

degree centrality result

EXAMPLE : FINDING INFLUENCER

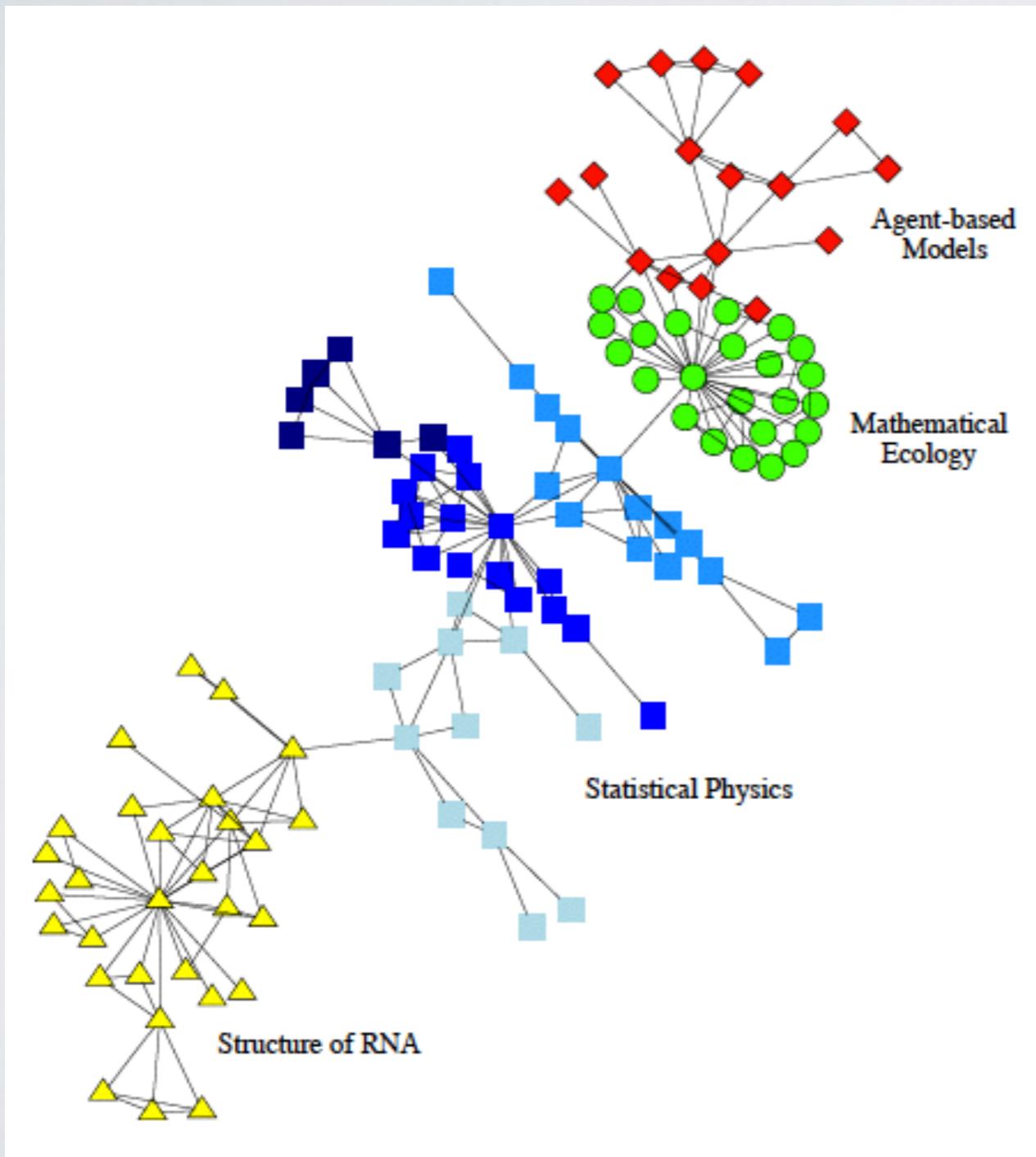


Geodesics	
length	#
Group Size	37
Potential Ties	1332
Actual Ties	170
Density	13%
4	558
5	136
6	0



Degrees		Betweenness		Closeness	
0.417	Mohamed Atta	0.334	Nawaf Alhazmi	0.571	Mohamed Atta
0.389	Marwan Al-Shehhi	0.318	Mohamed Atta	0.537	Nawaf Alhazmi
0.278	Hani Hanjour	0.227	Hani Hanjour	0.507	Hani Hanjour
0.278	Nawaf Alhazmi	0.158	Marwan Al-Shehhi	0.500	Marwan Al-Shehhi
0.278	Ziad Jarrah	0.116	Saeed Alghamdi*	0.480	Ziad Jarrah
0.222	Ramzi Bin al-Shibh	0.081	Hamza Alghamdi	0.429	Mustafa al-Hisawi

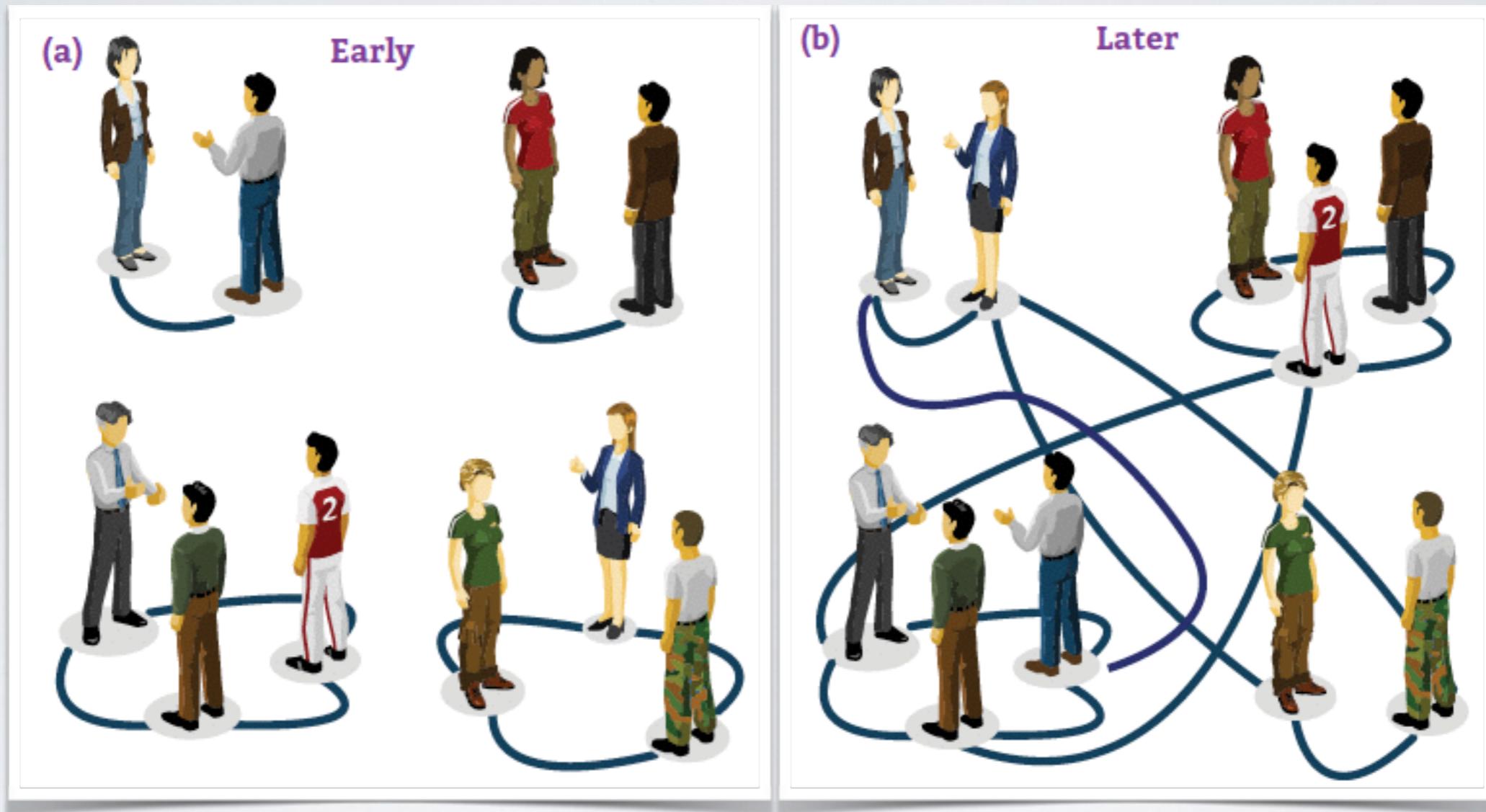
EXAMPLE : FINDING COMMUNITY



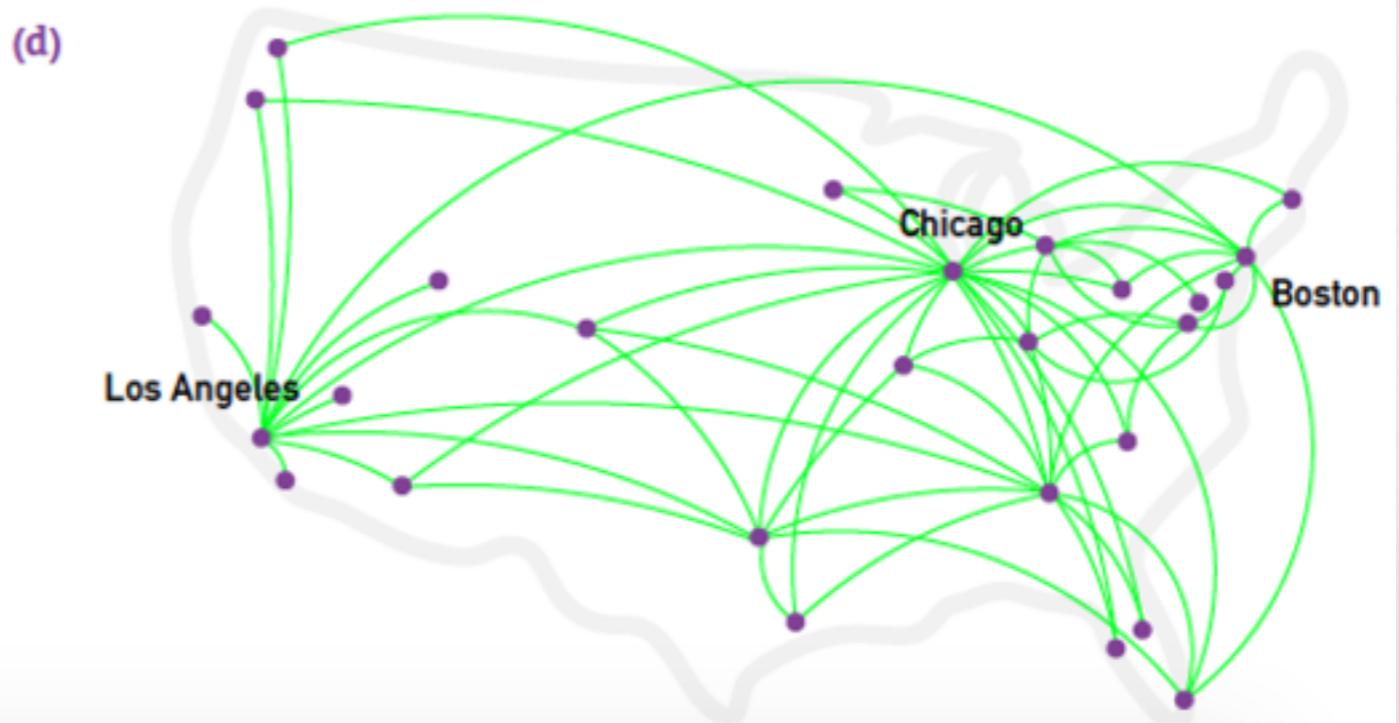
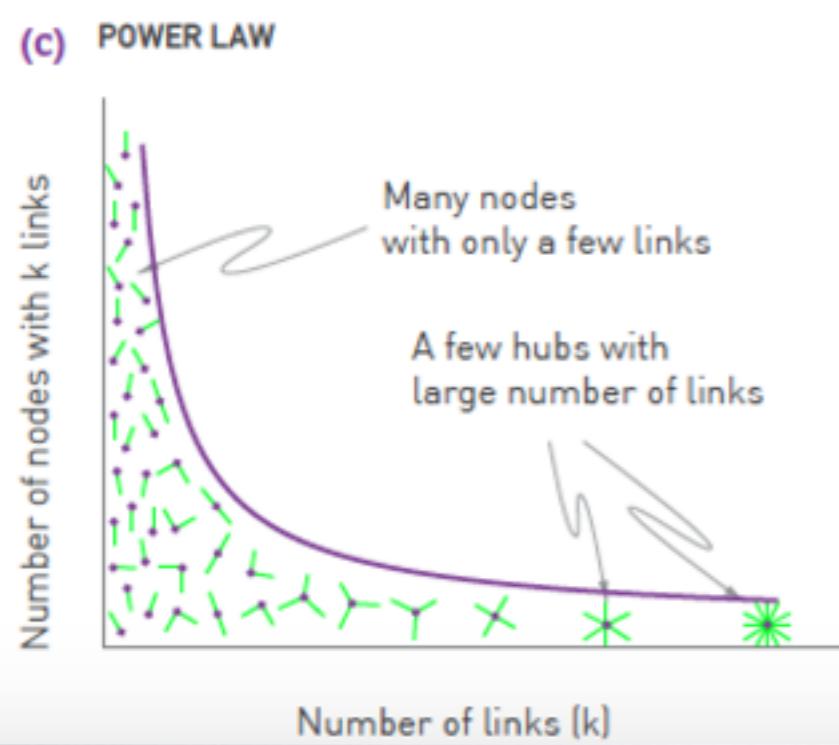
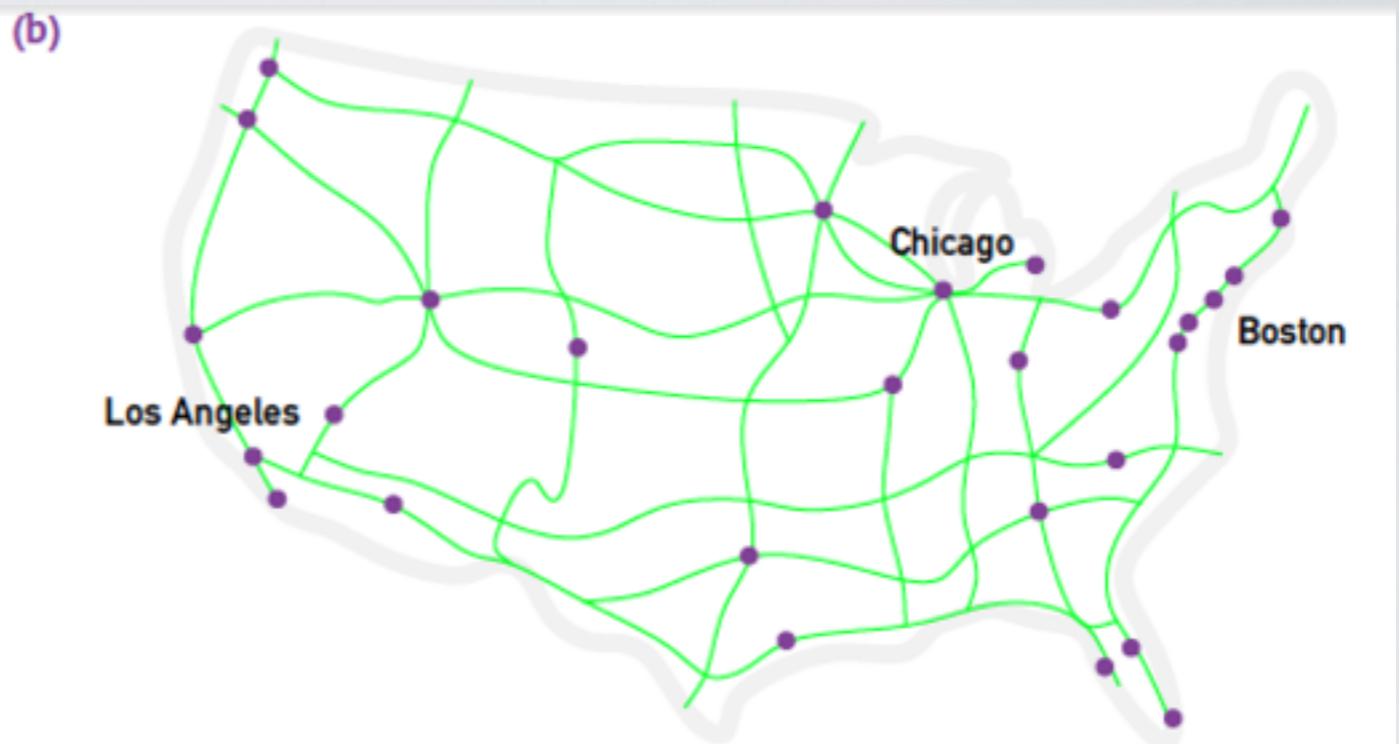
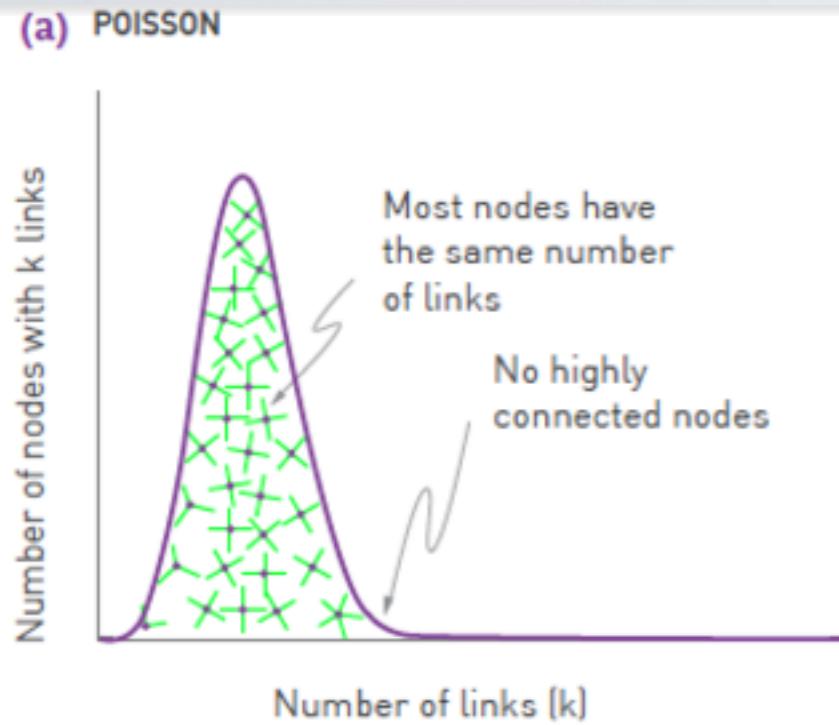
collaboration network of scientist at Santa Fe Institut (Girvan & Newman)
271 scientist (vertices) / 118 nodes from largest component
edge = scientist coauthor one of more publications

Komunitas : kumpulan titik titik dimana jumlah hubungan internal antar titik lebih besar dari pada jumlah hubungan dengan titik eksternal

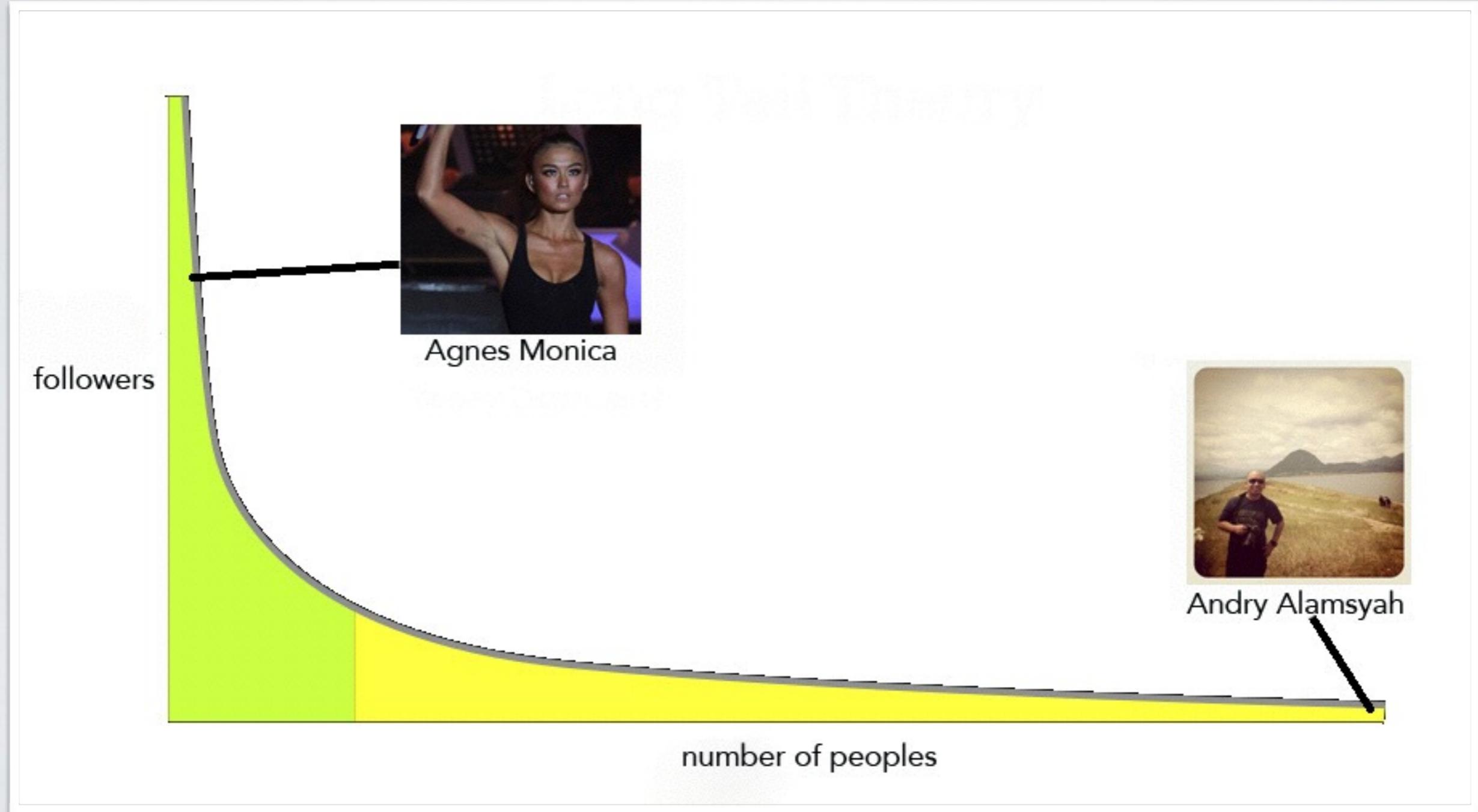
SOCIAL NETWORK FORMATIONS



SOCIAL NETWORK CHARACTERISTICS

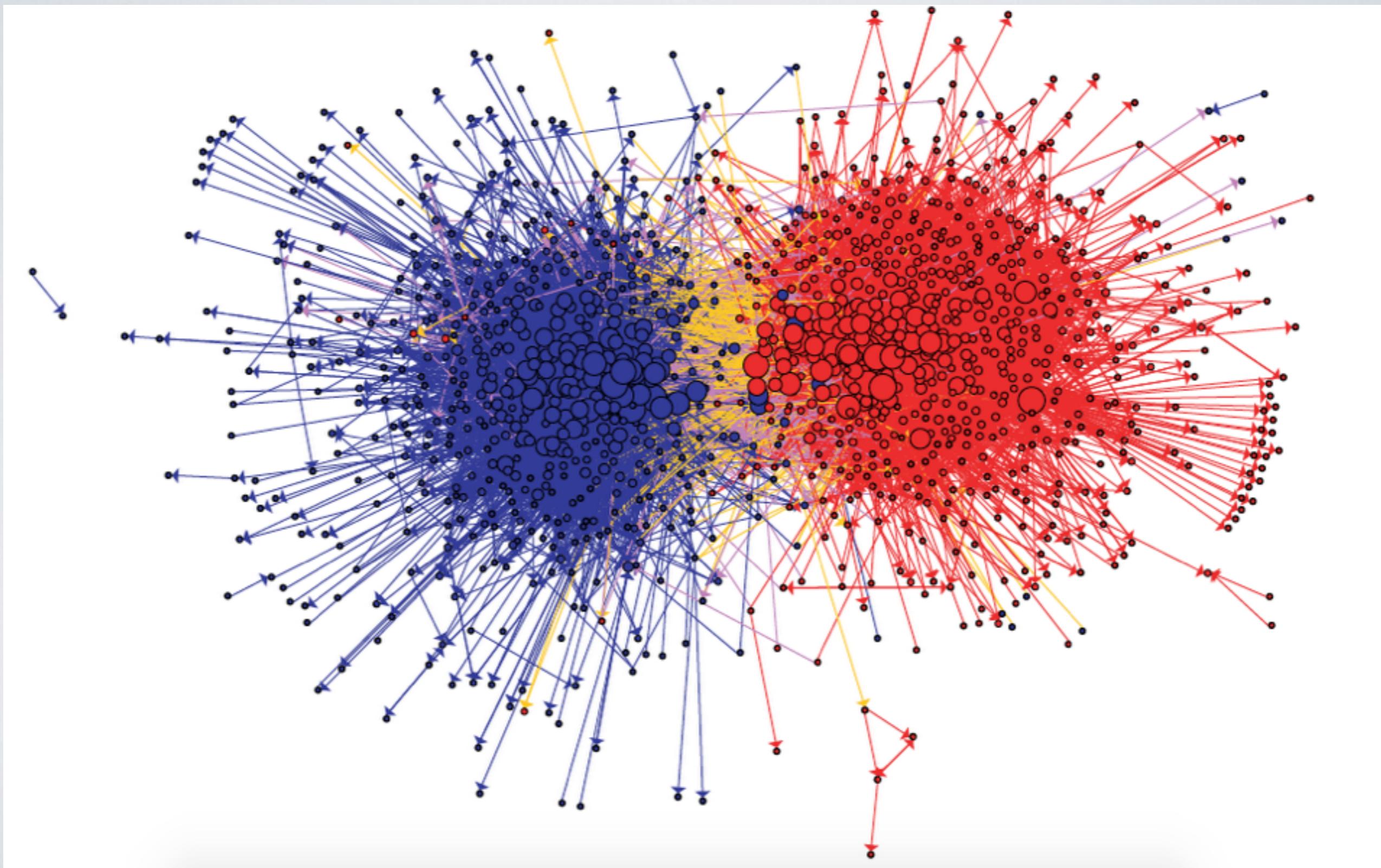


SOCIAL NETWORK CHARACTERISTICS



power law distributions

SOCIAL NETWORK CHARACTERISTICS

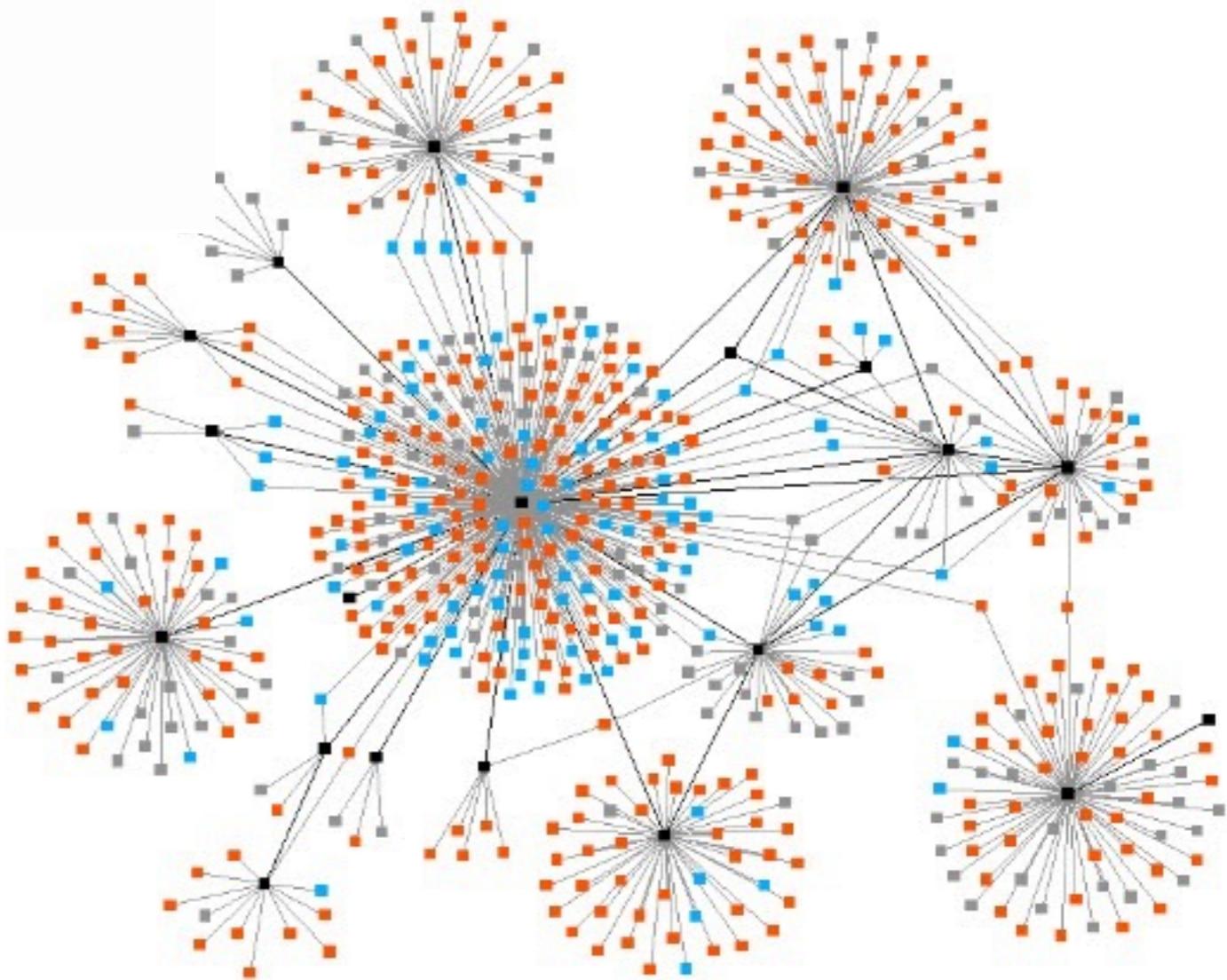


SOCIAL NETWORK CHARACTERISTICS

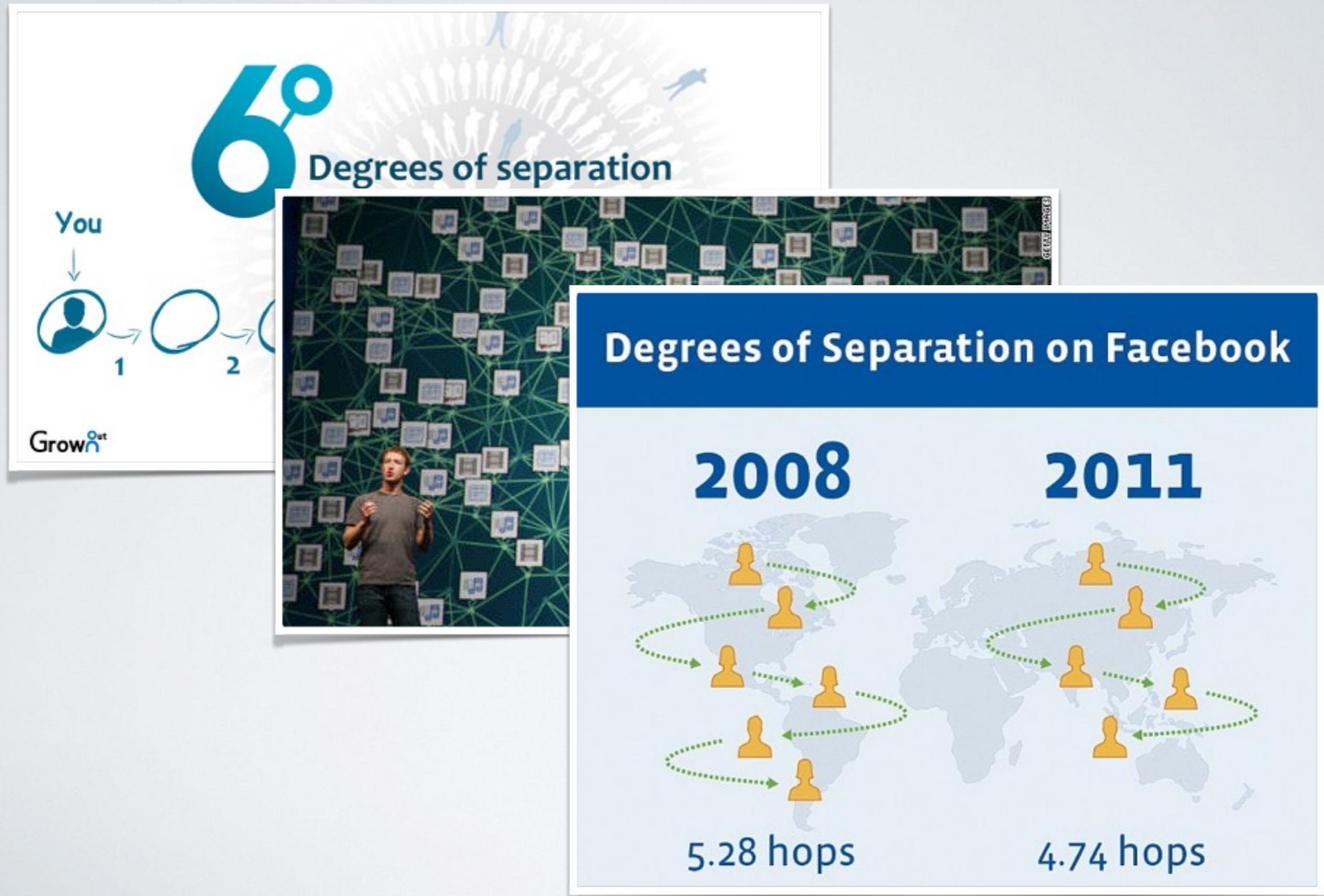


growth & preferential attachment

small world



SOCIAL NETWORK CHARACTERISTICS

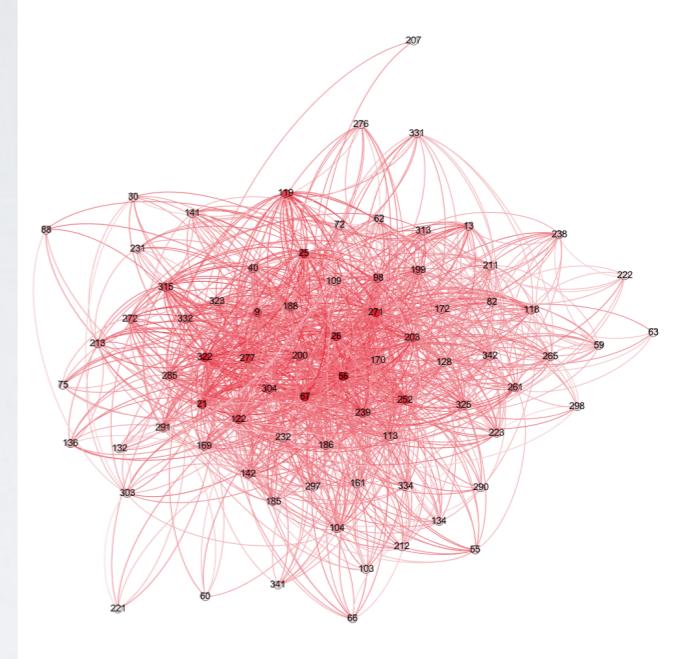
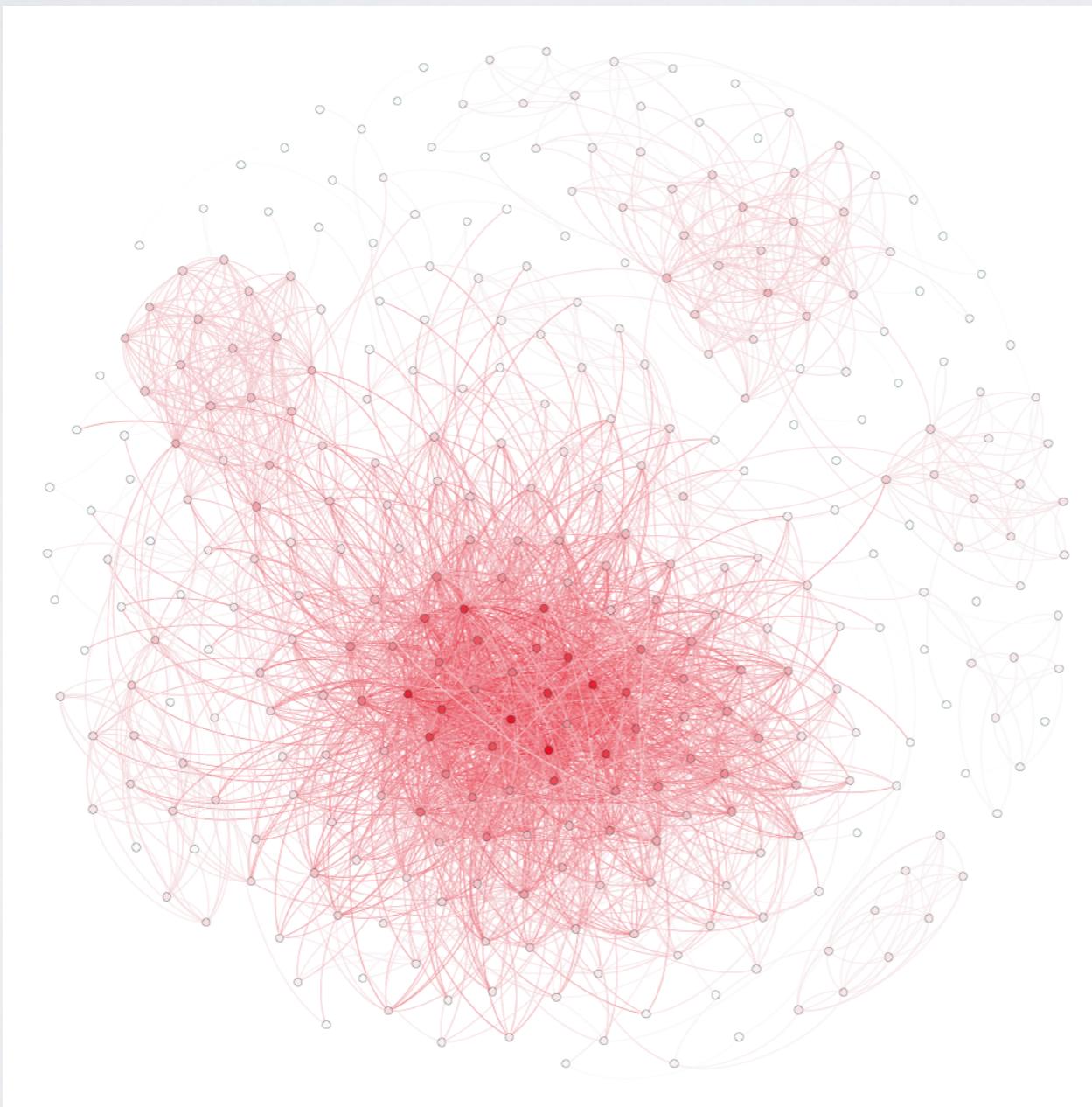


SOCIAL NETWORK INTERPRETATIONS

Case	Questions	SNA Tools
Leader Selection	Who is the central in the trust and respect network?	Degree Centrality
Ranks	How do we rank our top performer individuals in the organization?	Eigenvector Centrality, Pageranks
Task Force Selection	How do we put together a team that maximally connected through out the	Closeness Centrality
Mergers and Acquisition	How to merge separate cultures / networks?	Homophilly, Reciprocity, Mutuality, Transitivity
Competitive Advantages	What is the missing links between supply and demand?	Structural Holes
Advertising Attachment	How strong the impact of our advertisement effort?	Tie Strength
Market Segmentation	How segmented our market is?	Clustering Coefficient, Clique, Cohesive, Modularity
Information Dissemination	How is the information / knowledge spreading?	Random Walks, Hits Algorithm
Strength out the organization	How to increase redundancy and interconnectedness?	Bridge
Dynamics of Organization	How dynamics our organization is?	Temporal Networks

PUBLICATIONS (I)

Financial Fraud Detection using Social Network Analysis (eii-Forum, 2013)



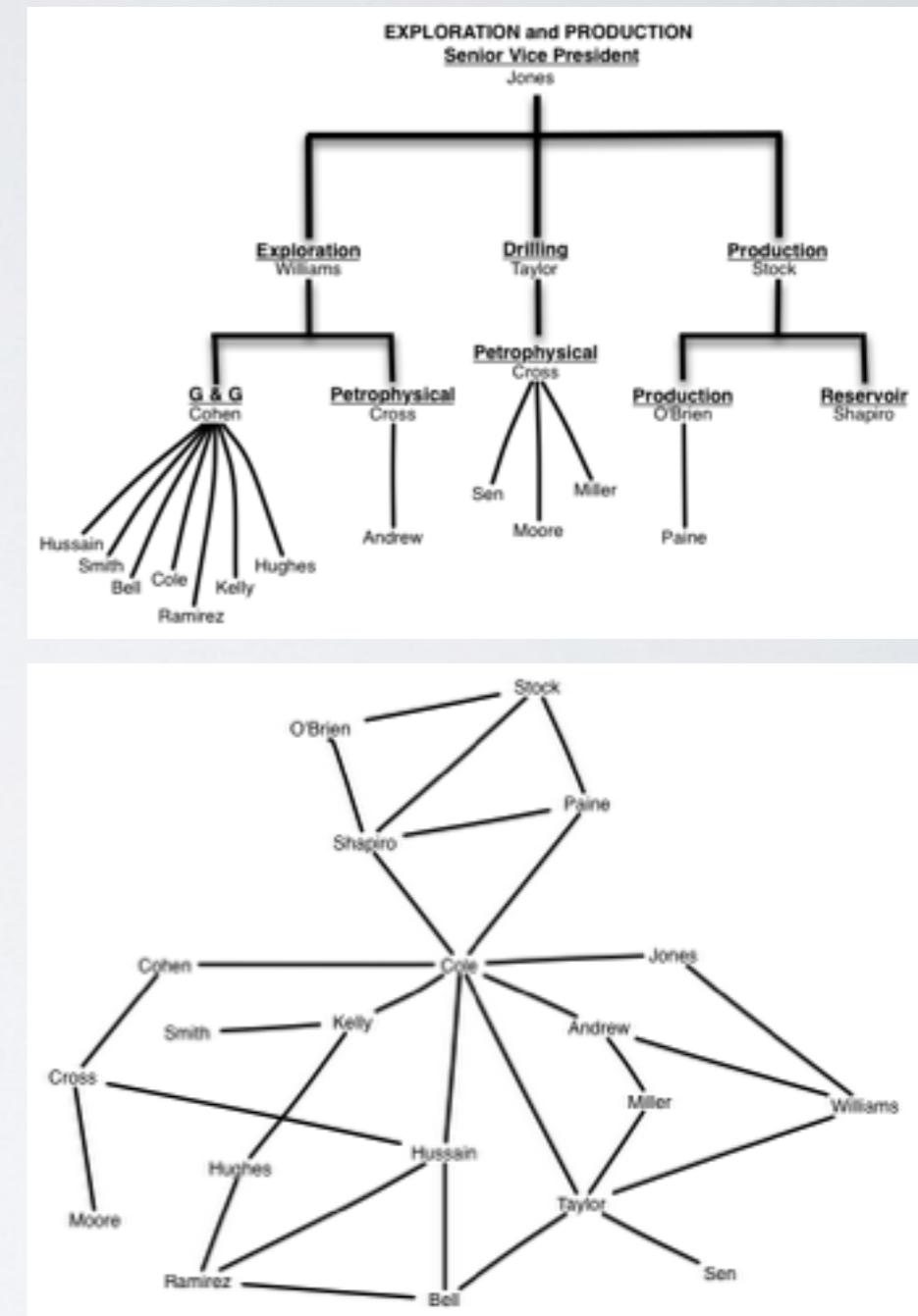
Tabel 3. komunitas, persentase dan node dengan degree terbesar pada komunitas.

Komunitas	Persentase (%)	Node dengan Degree Terbesar (a)	Degree Node (a)
2	31,23	56	77
7	13,51	312	25
3	12,61	175	16
1	12,31	271	72
0	10,21	53	30
4	8,11	119	61
5	3,9	320	20
8	3	275	9
10	2,4	89	7
9	0,9	179	2
12	0,6	233	1
11	0,6	282	1
6	0,6	33	1

PUBLICATIONS (2)

The Role of Social Network Analysis in Knowledge Management (Jurnal Manajemen Indonesia, 2013)

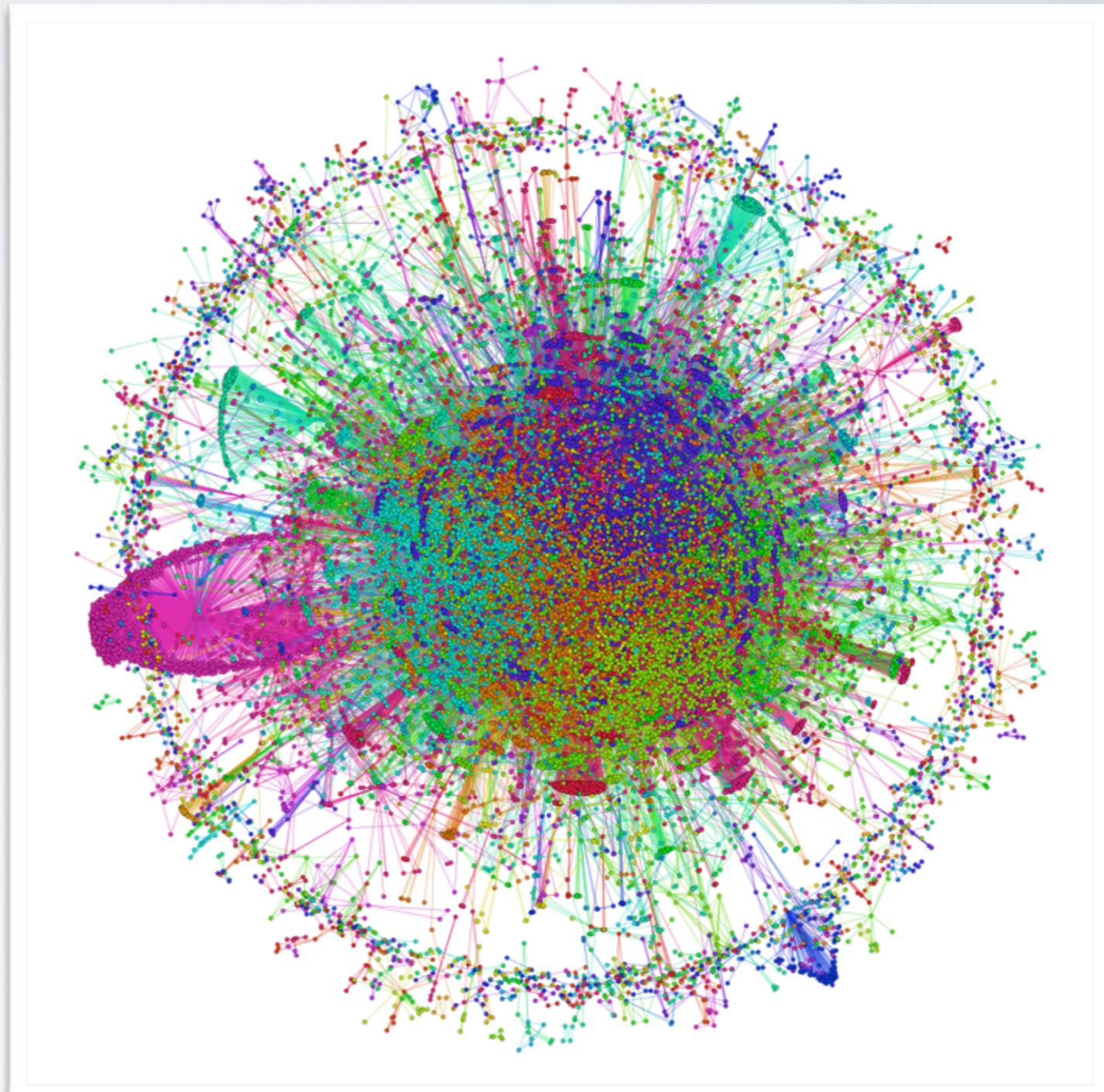
Case	Questions	SNA Tools
<i>Leader Selection</i>	<i>Who is the central in the trust and respect network?</i>	<i>Degree Centrality</i>
<i>Ranks</i>	<i>How do we rank our top performer individuals in the organization?</i>	<i>Eigenvector Centrality, Pageranks</i>
<i>Task Force Selection</i>	<i>How do we put together a team that maximally connected through out the organization?</i>	<i>Closeness Centrality</i>
<i>Mergers and Acquisition</i>	<i>How to merge separate cultures / networks?</i>	<i>Homophilly, Reciprocity, Mutuality, Transitivity</i>
<i>Competitive Advantages</i>	<i>What is the missing links between supply and demand?</i>	<i>Structural Holes</i>
<i>Advertising Attachment</i>	<i>How strong the impact of our advertisement effort?</i>	<i>Tie Strength</i>
<i>Market Segmentation</i>	<i>How segmented our market is?</i>	<i>Clustering Coefficient, Clique, Cohesive</i>
<i>Information Dissemination</i>	<i>How is the information / knowledge spreading?</i>	<i>Random Walks, Hits Algorithm</i>
<i>Strength out the organization</i>	<i>How to increase redundancy and interconnectedness?</i>	<i>Bridge</i>
<i>Dynamics of Organization</i>	<i>How dynamics our organization is?</i>	<i>Temporal Networks</i>



Vizualisation of hierarchical structure organization and knowledge flow of informal organization (Alamsyah, 2013)

PUBLICATIONS (3)

Effective Knowledge Management using Big Data and Social Network Analysis (ISCLO, 2013)



PUBLICATIONS (4)

Social Network Modelling Approach for Brand Awareness (ICOICT, 2014)

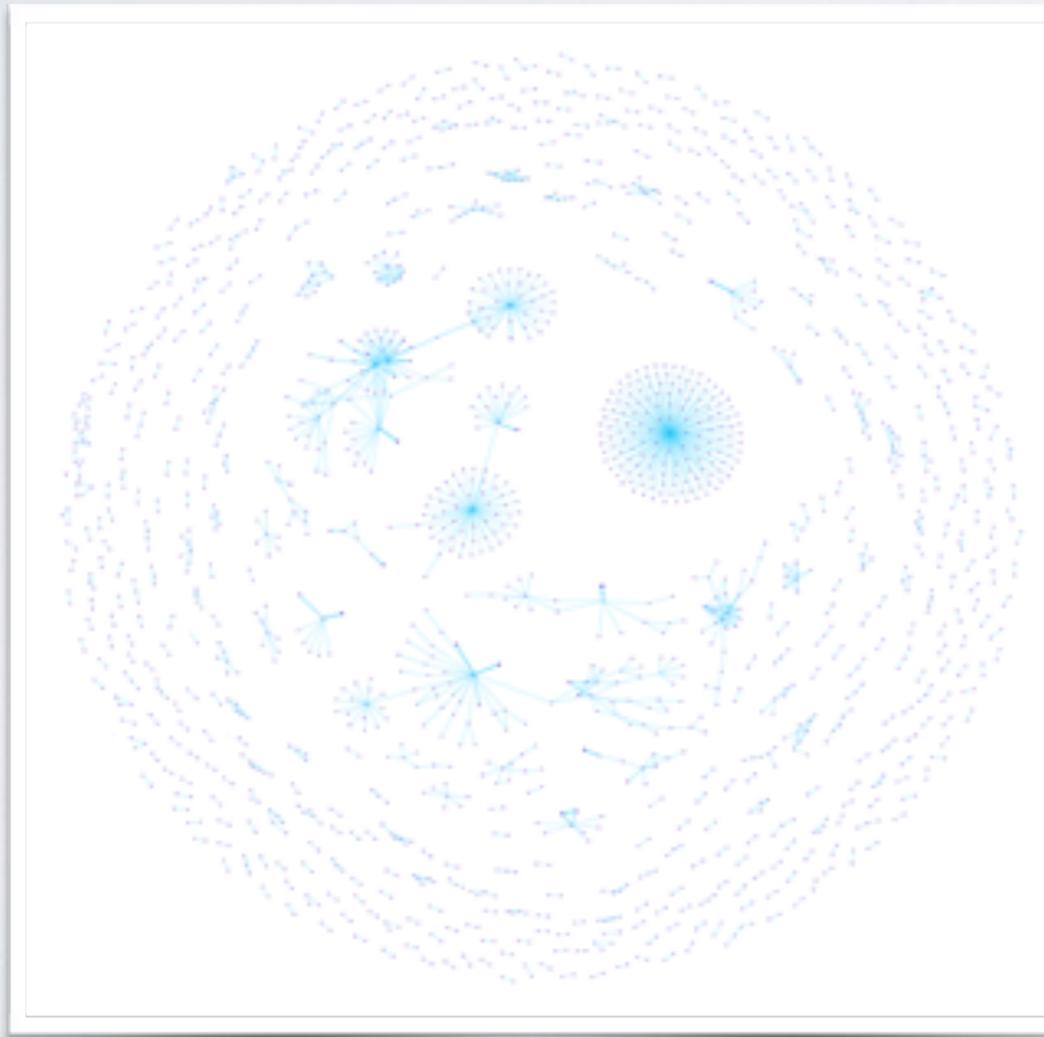


Figure 5. The graph model of twitter keyword “Telkom”

TABLE III. CENTRALITY RANK

No	Account	DC (value/rank)	BC (value/rank)	Community
1	@ouyalresearch	152 / 1	0.006 / 1	7
2	@detikcom	55 / 2	0.008 / 8	9
3	@infobdg	51 / 3	0.003 / 2	0
4	@EndankSoeka mti	43 / 4	0.002 / 3	1
5	@SID_Merch	42 / 5	0.002 / 4	1

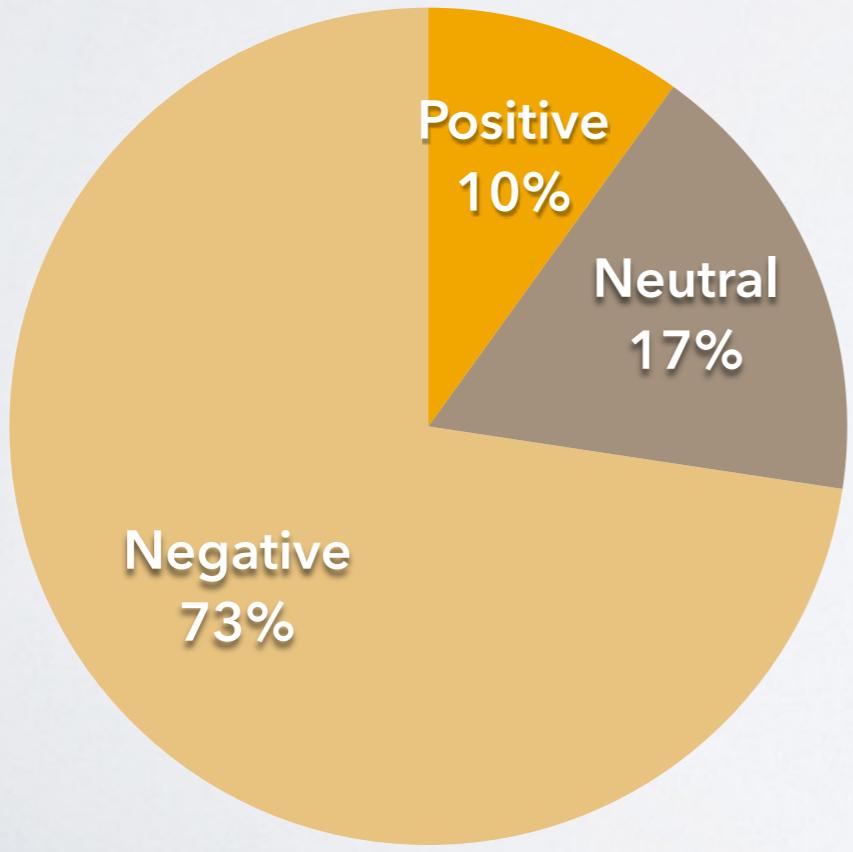
TABLE II. COMMUNITIES RANK

No	Community Number	Number of Node	Size (Percentage)
1	7	153	7.55
2	1	133	6.57
3	0	79	3.90
4	9	78	3.85
5	41	55	2.71
6	62	37	1.82
7	72	34	1.67
8	35	28	1.38

PUBLICATIONS (5)

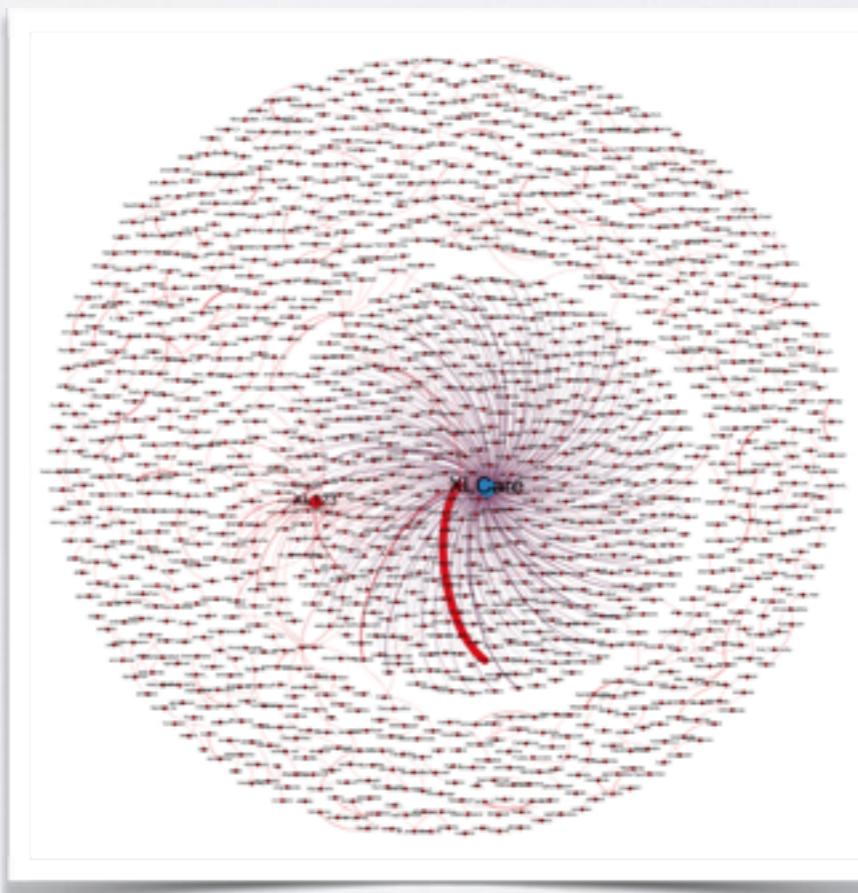
Topic : Social Customer Relationship Management based on Social Conversational Data (ISCLO, 2014)

SENTIMENT ANALYSIS



● Positive ● Neutral ● Negative

source	target
_Ad14	infobdg
AhmadAr	EndankSoekarmi
_AhmadR	KamtilBandung
_detikdotcom	detikdotcom
_dumb	KamtilStone
_Enanana	edofams
HTX23	indrawanthoni_
HTX23	indrawanthoni
_nd	onlydeha
Didi	putriel
SdxEka	mazaim22
transforming conversation data into social network	
Didi	pustridity
ouwresearch	
ITB_telkom	rachinien
ITB_telkom	rachinien
255_YUCHI	detikcom
Siti_nurjannah	IndiPreneur_ID
ze_e5	anitapramesty
Aashnugrah	endiahi
Aashnugrah	endiahi



PUBLICATIONS (6)

Network Market Analysis using Large Scale Social Network Conversation in Indonesia's Fast Food Industry (ICOICT, 2015)

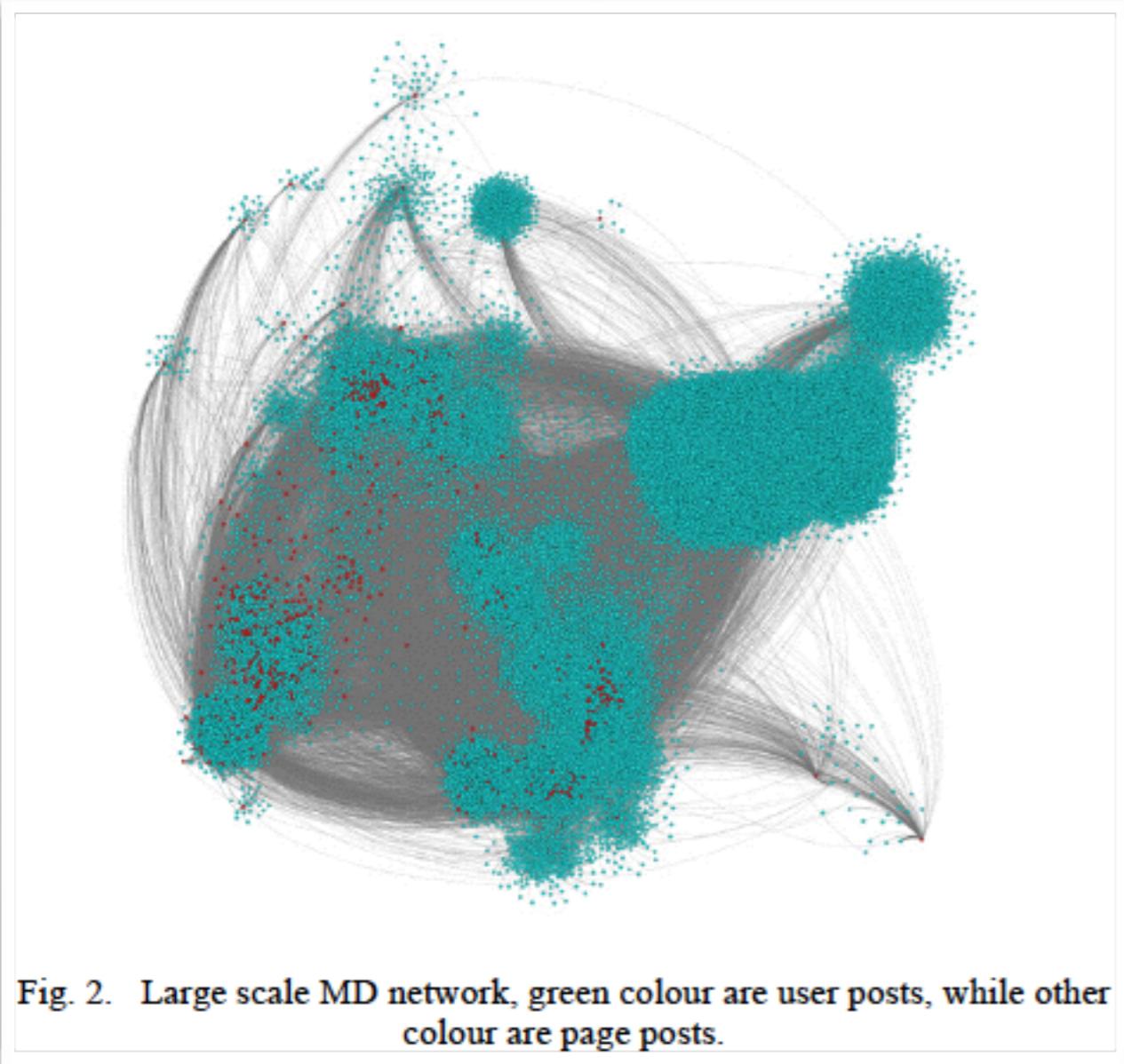


TABLE I. SNA METRICS COMPARISON

NO	METRIC	MD NETWORK	BK NETWORK
1	SIZE (NUMBER OF LIKES)	55,423,739	107,102
2	SIZE (BIPARTITE NETWORK)	438,568 NODES 909,748 EDGES	220,289 NODES 226,337 EDGES
3	AVERAGE DEGREE	2.418	0.116
4	DENSITY	0.00001097725	0.000003890513
5	AVERAGE PATH LENGTH	3.42506	3.245614
6	DIAMETER	8	8
7	MODULARITY	0.670 (13 COMMUNITIES)	0.527 (10 COMMUNITIES)

TABLE III. FURTHER INFORMATION ABOUT THE NETWORKS

NO	DESCRIPTION	MD NETWORK	BK NETWORK
1	POST TYPE	99.7% USER COMMENTS, THE REST ARE VIDEO, STATUS, PHOTO, LINK	99.55% USER COMMENTS, THE REST ARE VIDEO, STATUS, PHOTO, LINK
2	5 HIGHEST ENGAGEMENT VALUE	137609, 132372, 113299, 99668, 97973	129423, 25625, 25053, 21295, 17046
3	POST TYPE OF 5 HIGHEST ENGAGEMENT VALUE	PHOTO, STATUS, STATUS, LINK, PHOTO	ALL PHOTO
4	USER SEX	55.27% MALE, 41.50% FEMALE	55.68% MALE, 43.83% FEMALE

PUBLICATIONS (7)

Indonesian Music Fans Group Identification using Social Network Analysis in Kaskus Forum (ICOICT, 2015)

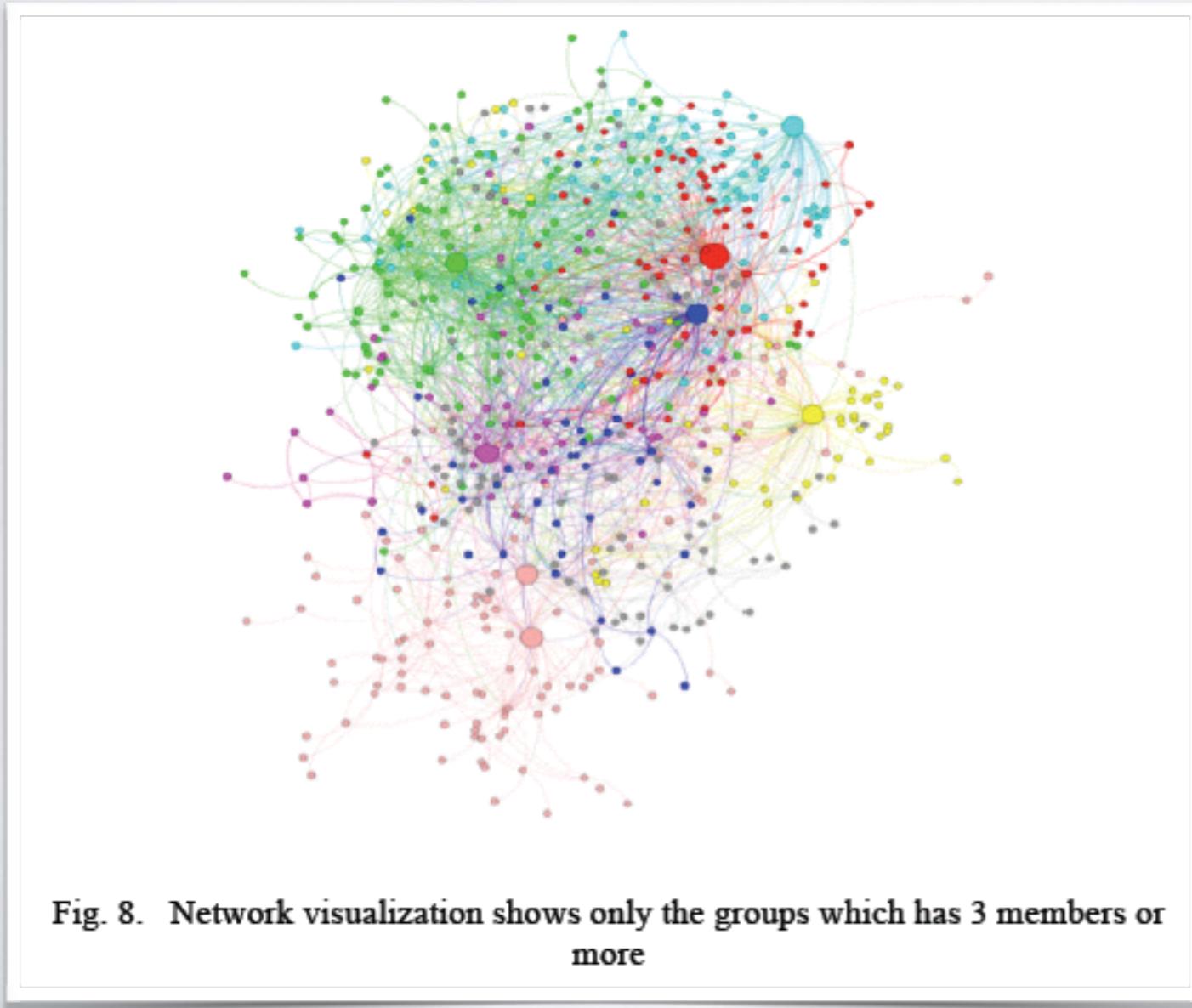
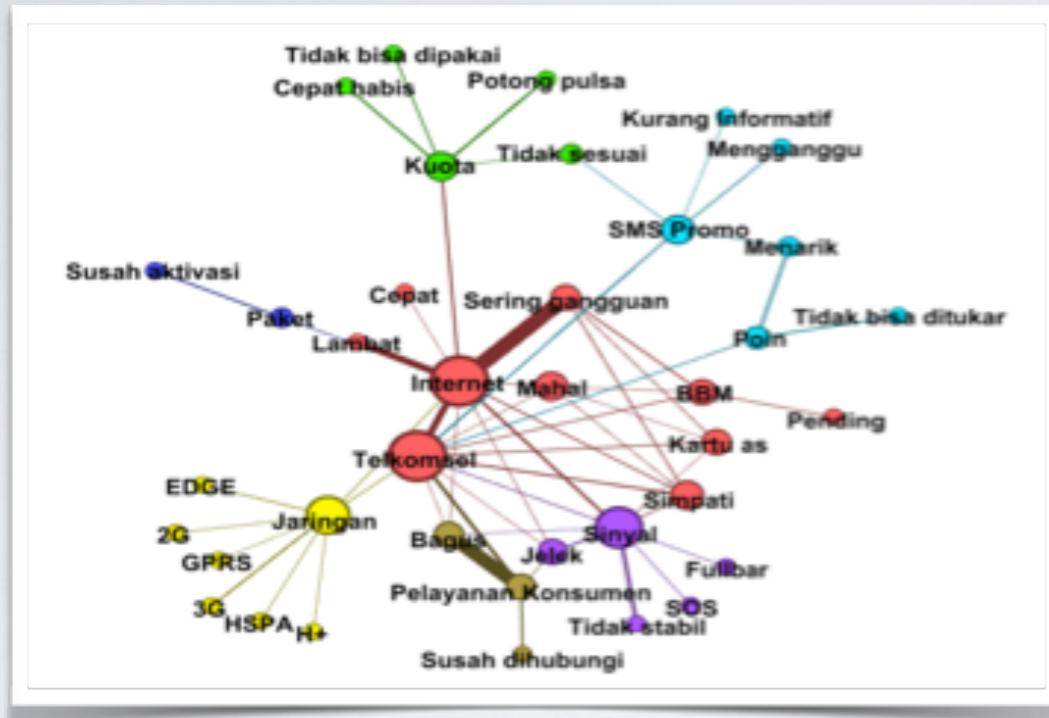


Fig. 8. Network visualization shows only the groups which has 3 members or more

TABLE III. SIZE DISTRIBUTION PERCENTAGE

Size	The Number of Groups	Percentage
123	1	9,34%
115	1	8,73%
78	1	5,92%
76	1	5,77%
56	1	4,25%
50	1	3,8%
49	1	3,72%
45	1	3,42%
40	1	3,04%
15	1	1,14%
3	1	0,23%
2	8	0,15% x 8
Total	19	50.56%

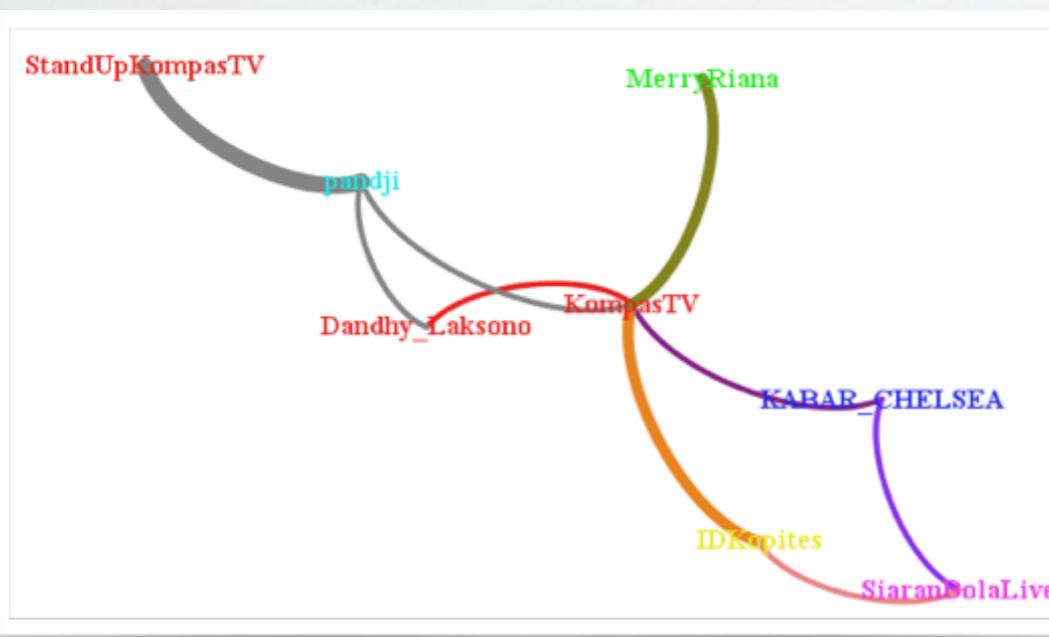
UNPUBLISHED WORK



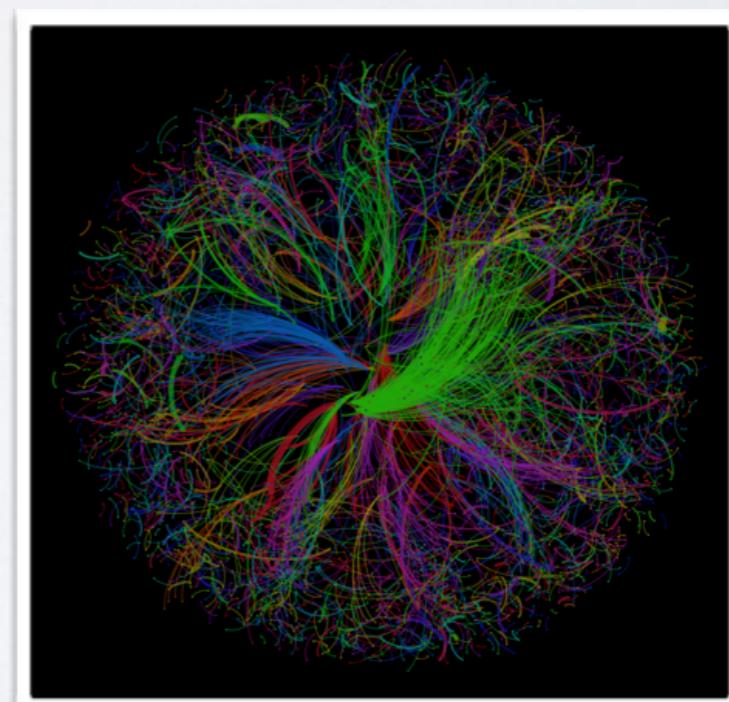
Social Network Analysis, Text Mining, Sentiment Analysis for CRM (Telkomsel)



Telkom University Academics Network (academia.edu)



Influencer Analysis for Alternative Marketing strategy of KompasTV



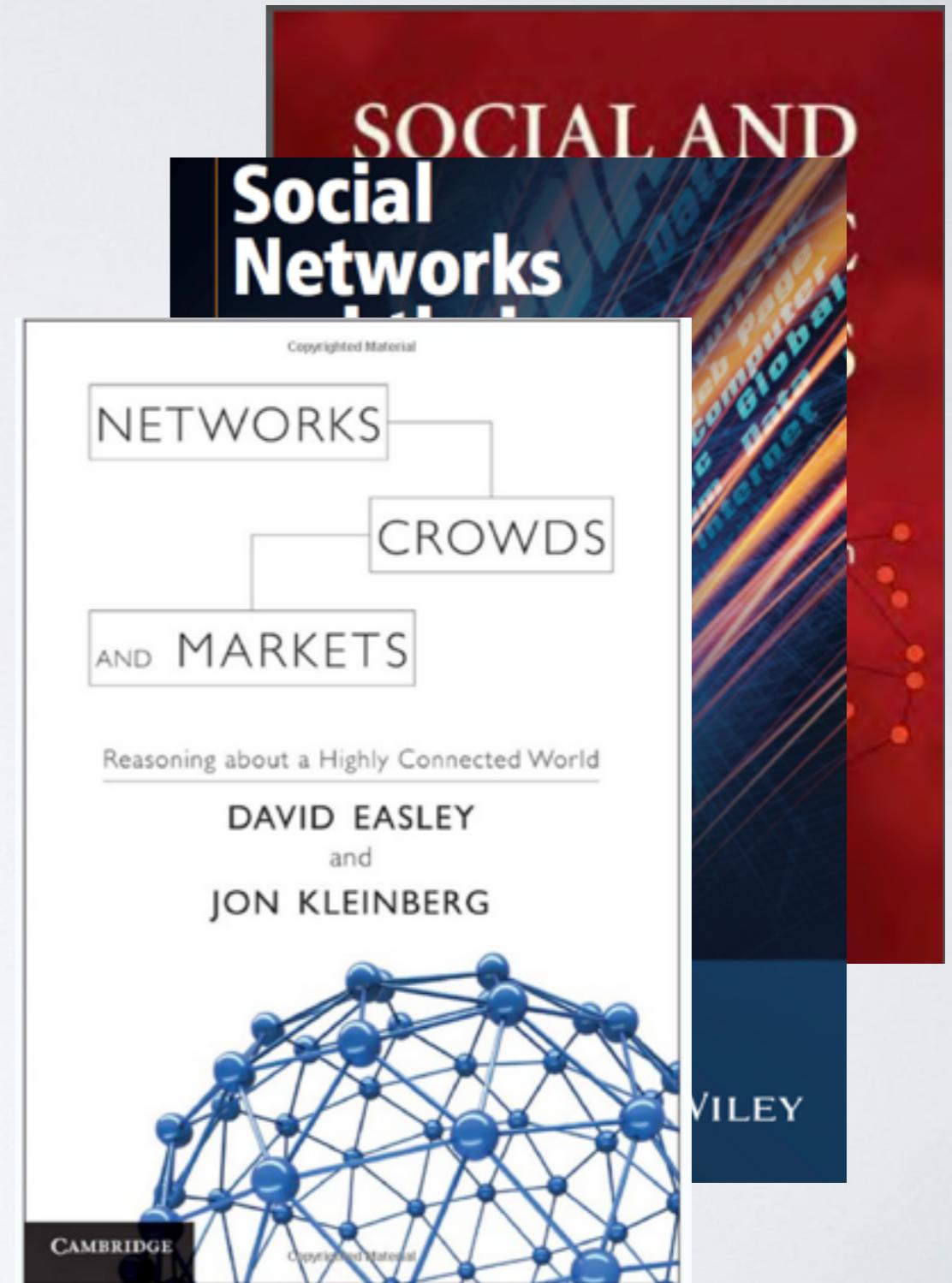
Word of Mouth Simulation Process on Hashtag #Android #BBM

ONGOING RESEARCH

- *Social Engagement Analysis* in online conversation of higher education community; use Facebook data
- *Top of Mind (Brand Awareness)* methods comparison using legacy approach and social media crawling data; verification research
- *Sentiment Analysis* using appraisal theory; mobile phone product comparison
- New *Metric Constructions* for engagement and dissemination information (example: marketing purpose)
- *Smart City Metric* based on social media
- Classification models for online SME selling and production effort
- Data mining modelling for online advertising and online tourism
- Indonesian investor and start up network model. Based on small world phenomenon
- Text network analysis comparison for telco industry
- Online conversation comparison of three biggest bank in Indonesia
- etc ...

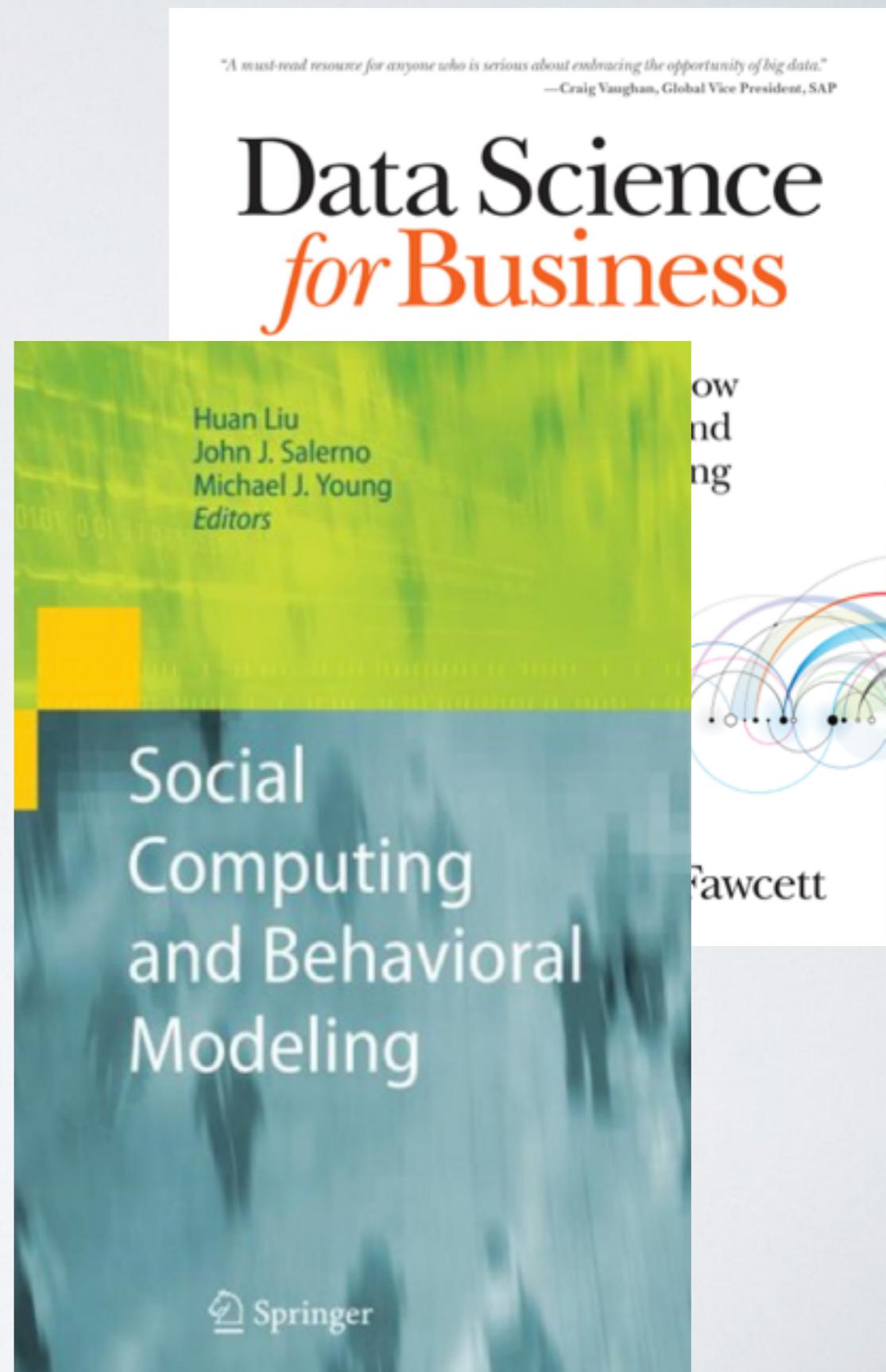
OPPORTUNITY

- dynamics network (size, movement, temporal, evolving network)
- network market predictions
- social network aggregate behaviour (game theory on traffic, market, auctions)
- bargaining and powers in networks
- social and economic network (strategic formation, diffusion, learning from networks, etc)
- social network and communication pattern in the context of smart city
- metric construction
- graph algorithm
- etc ..



CONCLUSIONS

- The availability of large volume dataset should encourage business and social science (as well as other sciences) to apply *Big Data* (SNA/DM/KDD/DS) tools
- Business and Social Science (as well as other sciences) need to use more *Big Data* tools, because of the capability to model, predict, simulate and optimise phenomenon
- There are many vacancy in Social Network (Graph Representation) research, such as large scale graph, graph reduction, graph algorithm
- *Big Data* approach are more and more utilised as standard of decision making in modern business



without Big Data, you are blind and deaf in the middle of a
freeway - Geoffrey Moore

TERIMA KASIH