

FUNDAMENTALS OF (ONLINE ((SOCIAL) MEDIA))) NETWORK ANALYSIS

LECTURE 1

Online Social Networks, Elements of Networks, Network Measures, Data Sources

WHO IS THIS GUY?

WHY ARE YOU HERE?

Who are you?

Why did you choose this course?

What are your expectations for this day?

THE PLAN

1. Online (Social) Media Network Fundamentals
2. Network fundamentals
3. *break*
4. Network Analysis Methods
5. Data Mining Possibilities and Difficulties

Afterwards:

Practical on Data Collection and Exploratory Analysis with Descriptive Statistics in Python

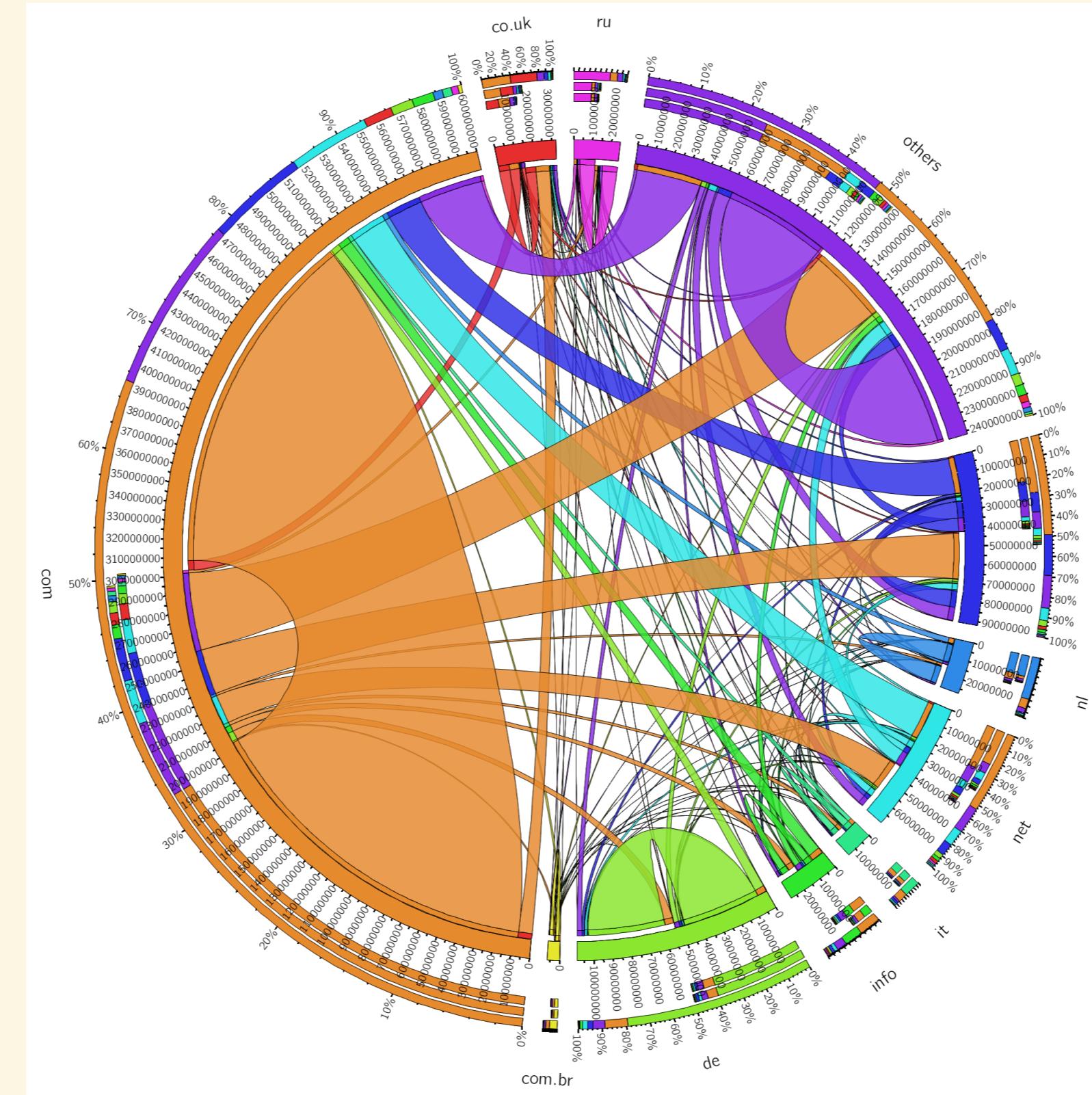
ONLINE ((SOCIAL) MEDIA) NETWORKS

NOT TECHNICAL INFRASTRUCTURE NETWORKS

The Internet: 1997 - 2021 Contextual



MOSTLY NOT HYPERLINK NETWORKS



Links between top level domains in 2012 ("Topology of the WDC Hyperlink Graph",
<http://km.aifb.kit.edu/sites/webdatacommons/hyperlinkgraph/topology.html>)

INFLUENCE (SOME/MOST?) INFORMATION DIFFUSION

#SYDNEYSIEGE VS #ILLRIDEWITHYOU VS \BREXIT PETITION

Münch, F. V. (2019). *Measuring the Networked Public – Exploring Network Science Methods for Large Scale Online Media Studies* [PhD thesis, Queensland University of Technology].
<https://doi.org/10.5204/thesis.eprints.125543>

#ILLRIDEWITHYOU

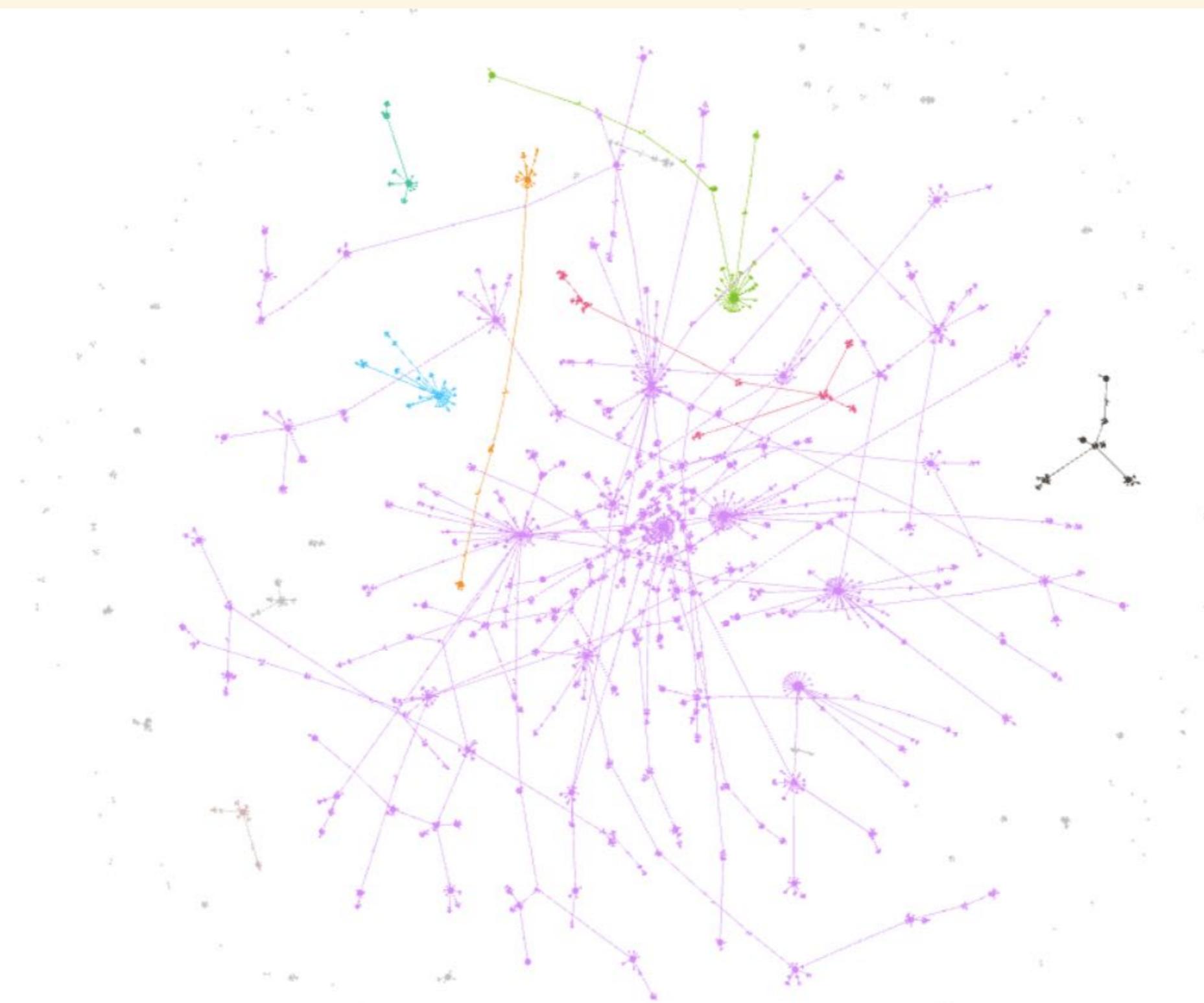
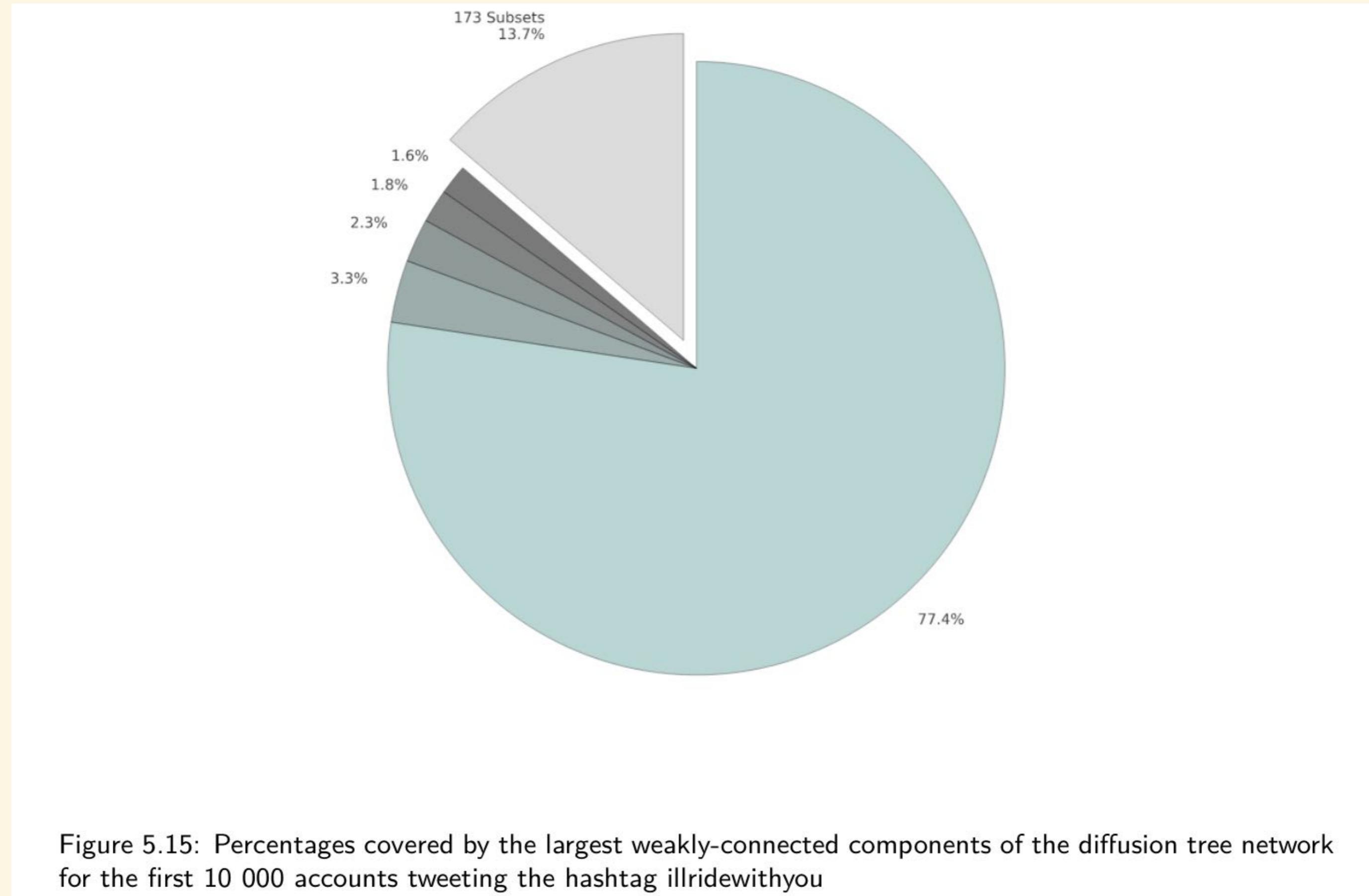
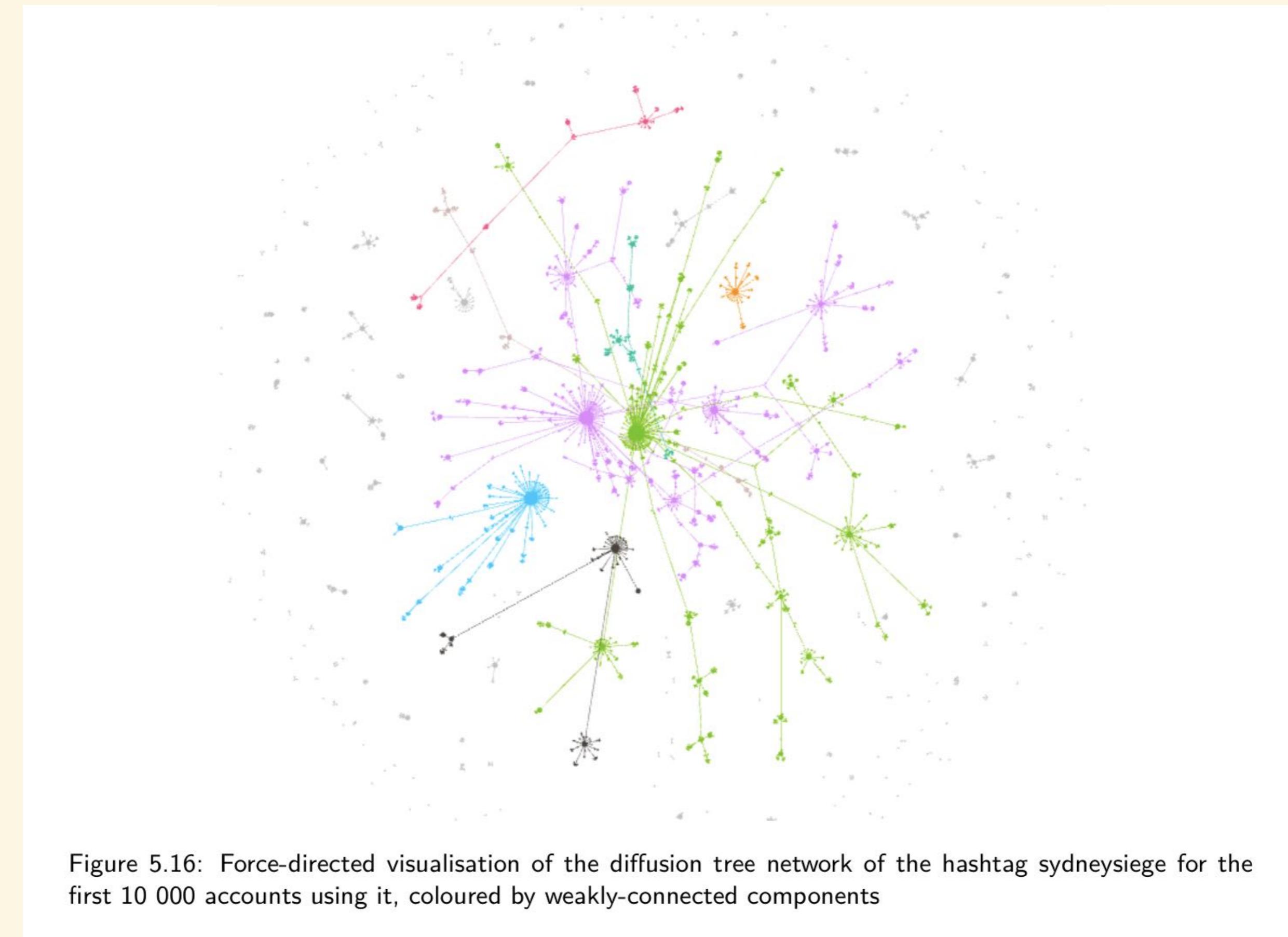
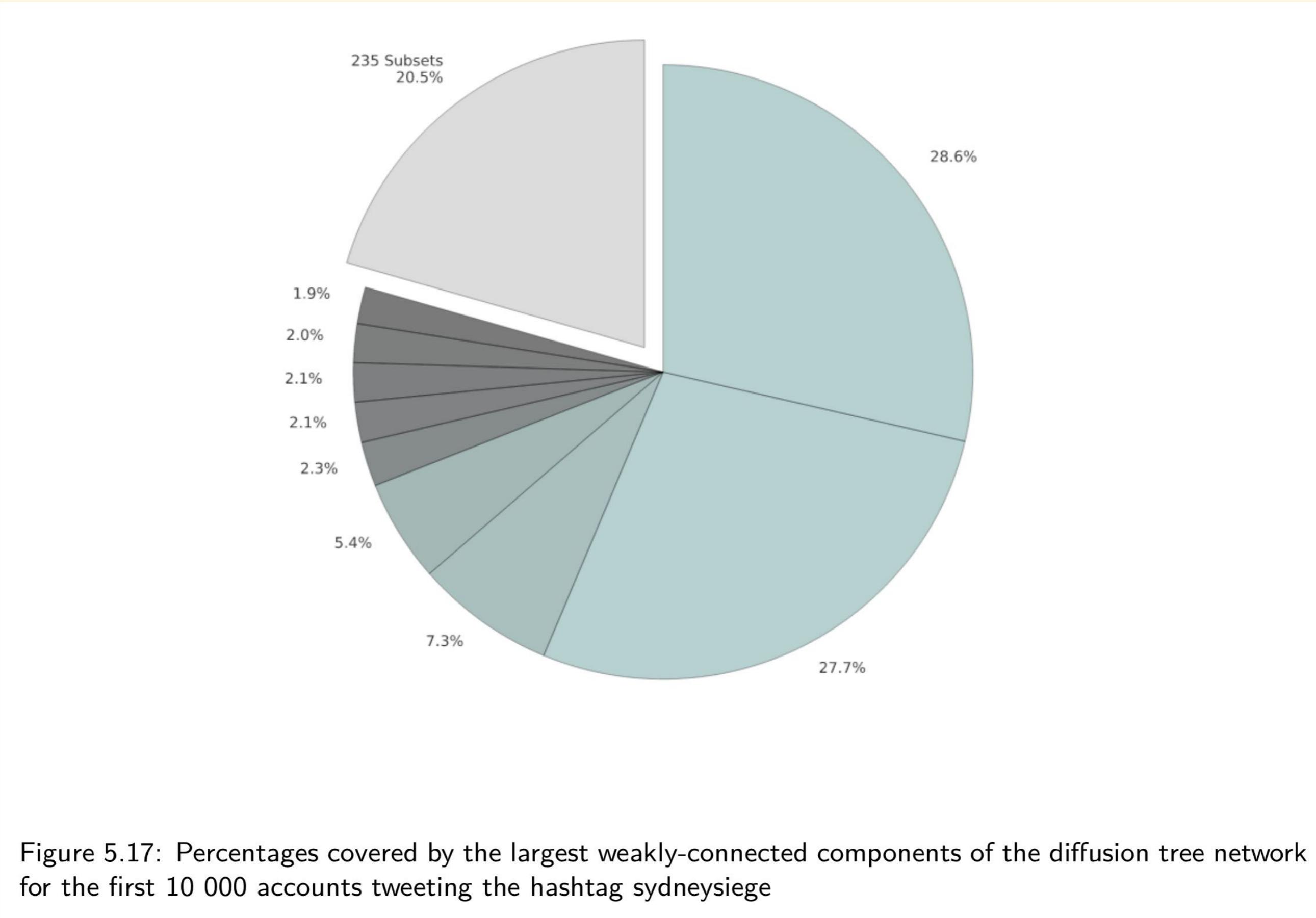


Figure 5.14: Force-directed visualisation of the diffusion tree network of the hashtag illridewithyou for the first 10 000 accounts using it, coloured by weakly-connected components

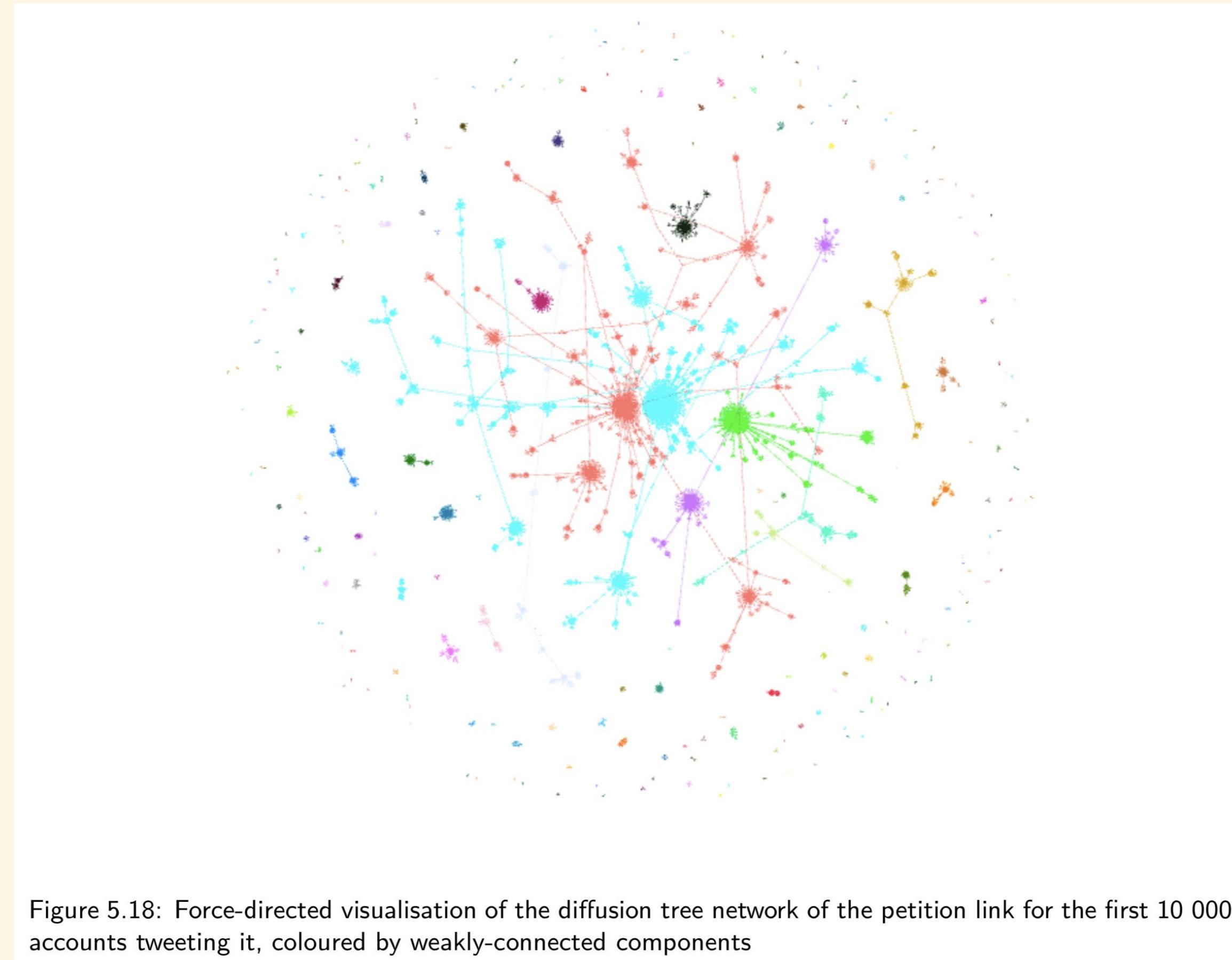


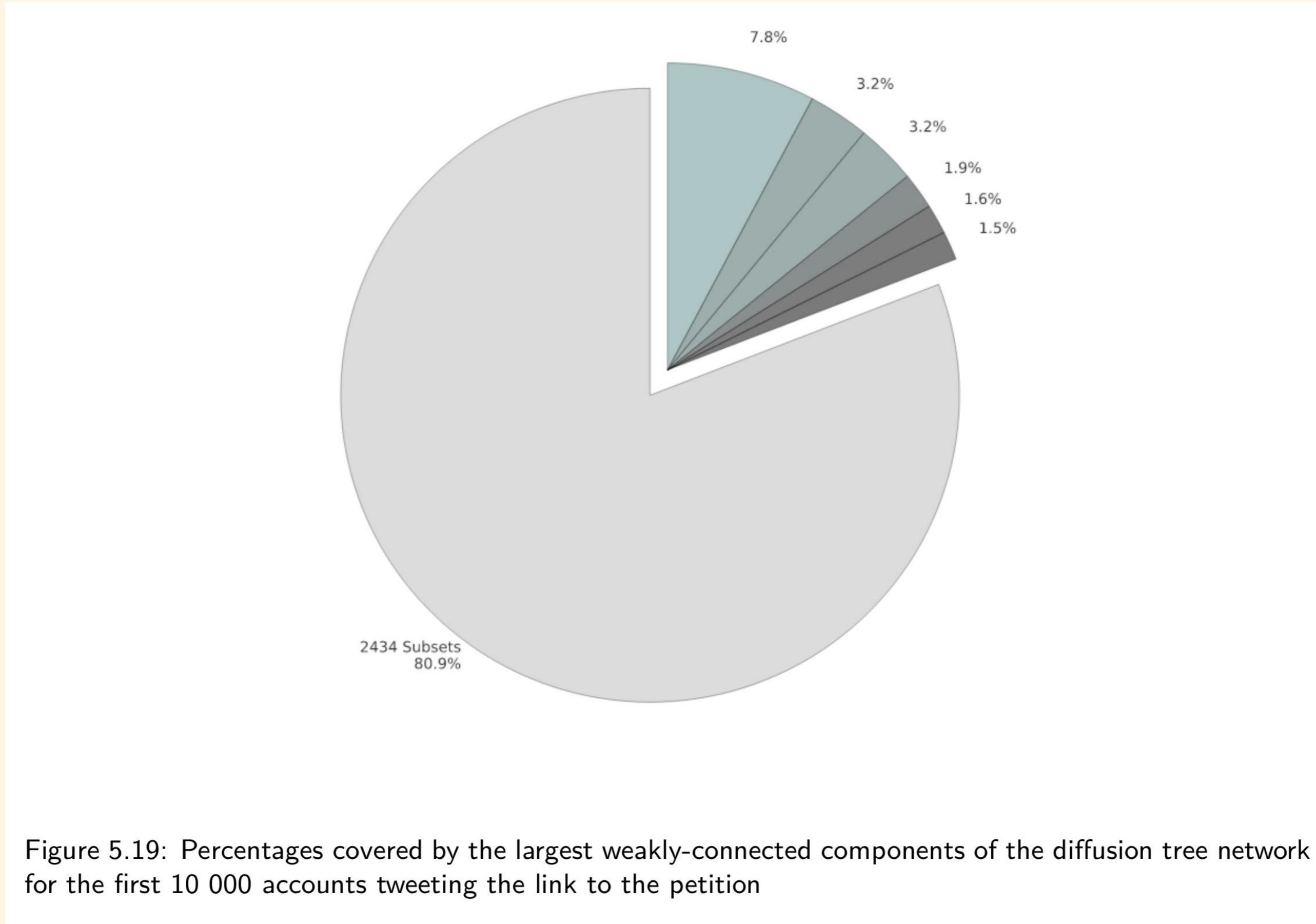
#SYDNEYSIEGE





ANTI-BREXIT PETITION





**THE NUMBER AND SIZE OF CONNECTED COMPONENTS INDICATES THE INFLUENCE OF THE NETWORK COMPARED TO
OUTSIDE SOURCES**

LEAD TO CLASSIFIABLE COMMUNICATION PATTERNS

Himelboim, I., Smith, M. A., Rainie, L., Shneiderman, B., & Espina, C. (2017). *Classifying Twitter topic-networks using social network analysis*. 1–38. <https://doi.org/10.1177/2056305117691545>

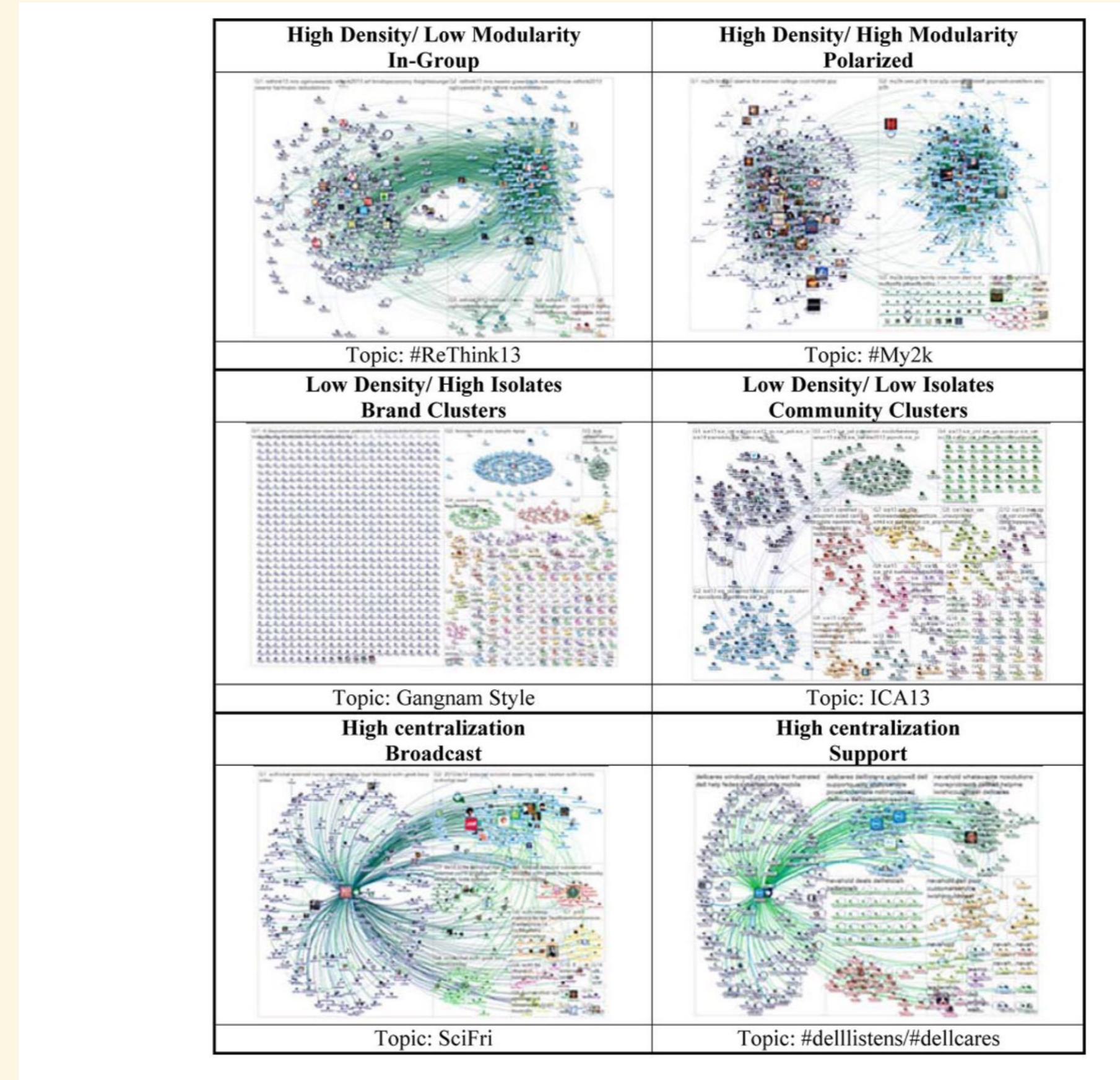
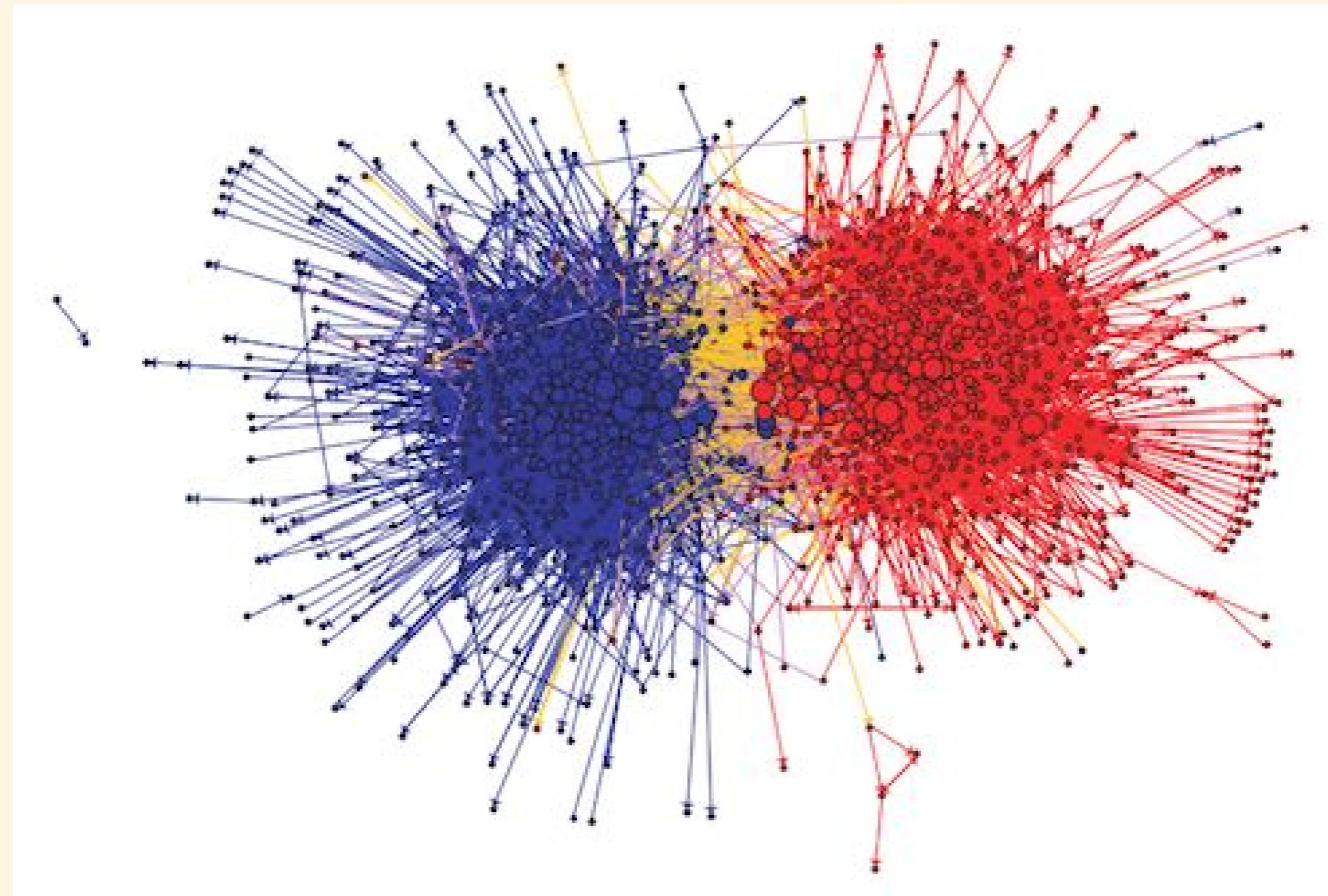


Figure 3. Network visualization by topic-network category.

EXAMPLE: POLARISATION

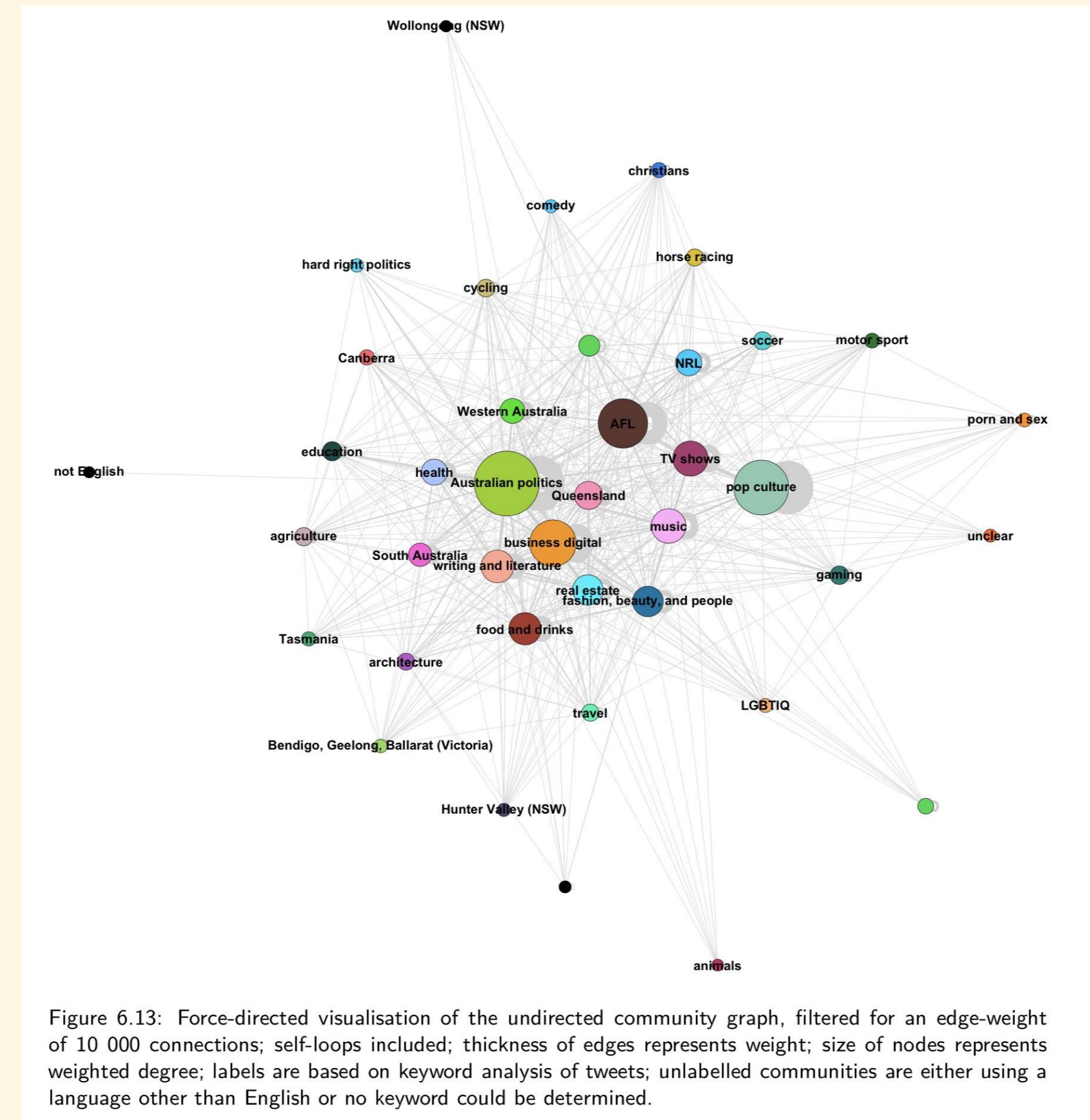


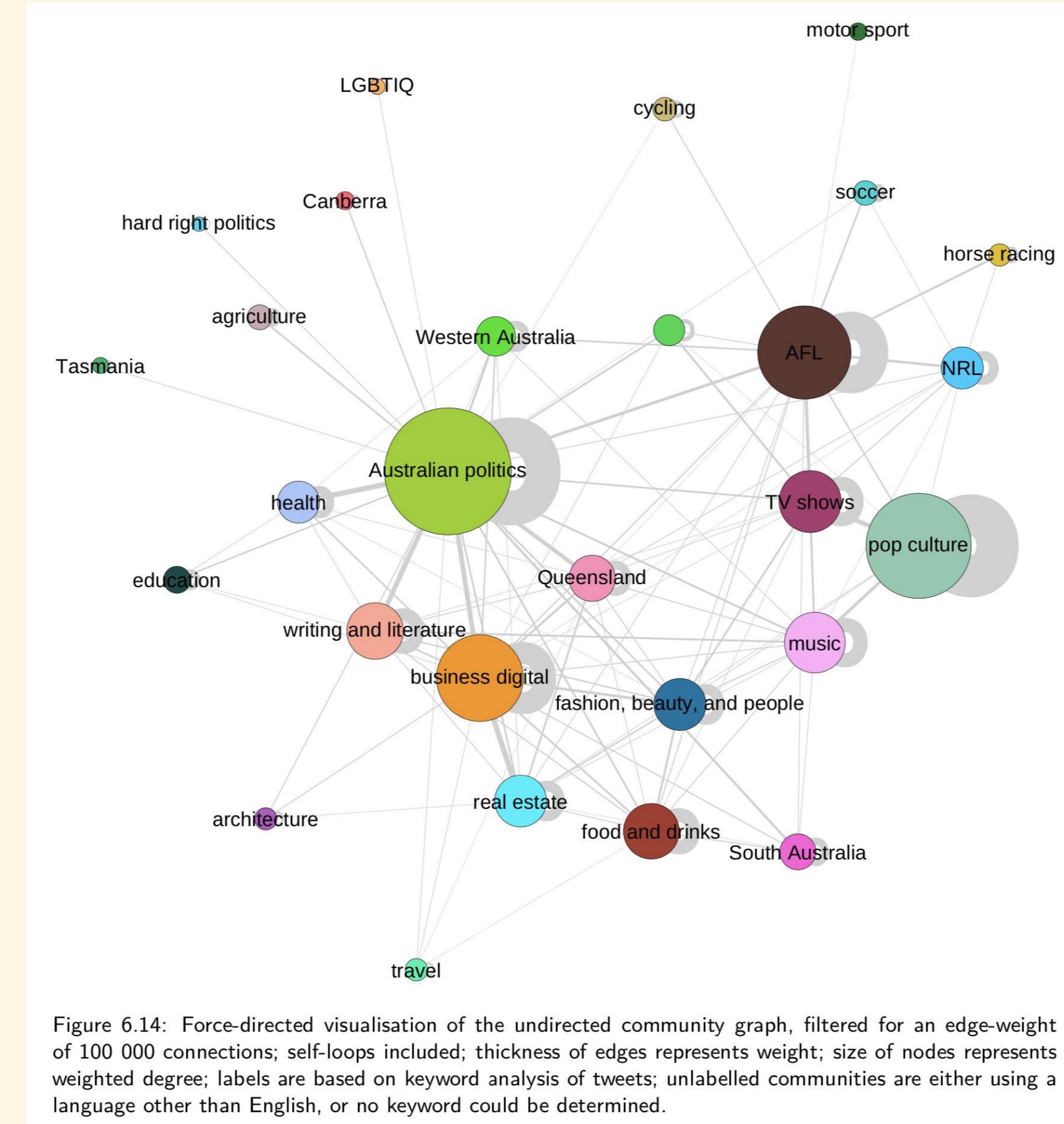
Lada A. Adamic and Natalie Glance. 2005. The political blogosphere and the 2004 U.S. election: divided they blog. In Proceedings of the 3rd international workshop on Link discovery (LinkKDD '05). Association for Computing Machinery, New York, NY, USA, 36–43. <https://doi.org/10.1145/1134271.1134277>

REFLECT LONG-TERM STRUCTURED SYSTEMS OF (PARTS OF) SOCIETY

- Bruns, A., Moon, B., Münch, F. V., & Sadkowsky, T. (2017). The Australian Twittersphere in 2016: Mapping the follower/followee network. *Social Media + Society*, 3(4).
<https://doi.org/10.1177/2056305117748162>
- Münch, F. V. (2019). *Measuring the Networked Public – Exploring Network Science Methods for Large Scale Online Media Studies* [PhD thesis, Queensland University of Technology].
<https://doi.org/10.5204/thesis.eprints.125543>
- Münch, F. V., & Rossi, L. (2020, October 5). A Tale of Two Twitters? Identifying Bridges Between Language Based Twitterspheres. *AoIR Selected Papers of Internet Research*.
<https://doi.org/10.5210/spir.v2020i0.11283>
- Münch, F. V., & Rossi, L. (2020). *Bootstrapping Follow Networks of Influential Twitter Accounts*. IC2S2. <https://vimeo.com/431470176>
- Münch, F. V., Thies, B., Puschmann, C., & Bruns, A. (2021). Walking Through Twitter: Sampling a Language-Based Follow Network of Influential Twitter Accounts. *Social Media + Society*.
<https://doi.org/10.1177/2056305120984475>

AUSTRALIAN TWITTERSPHERE





GERMAN TWITTERSPHERE

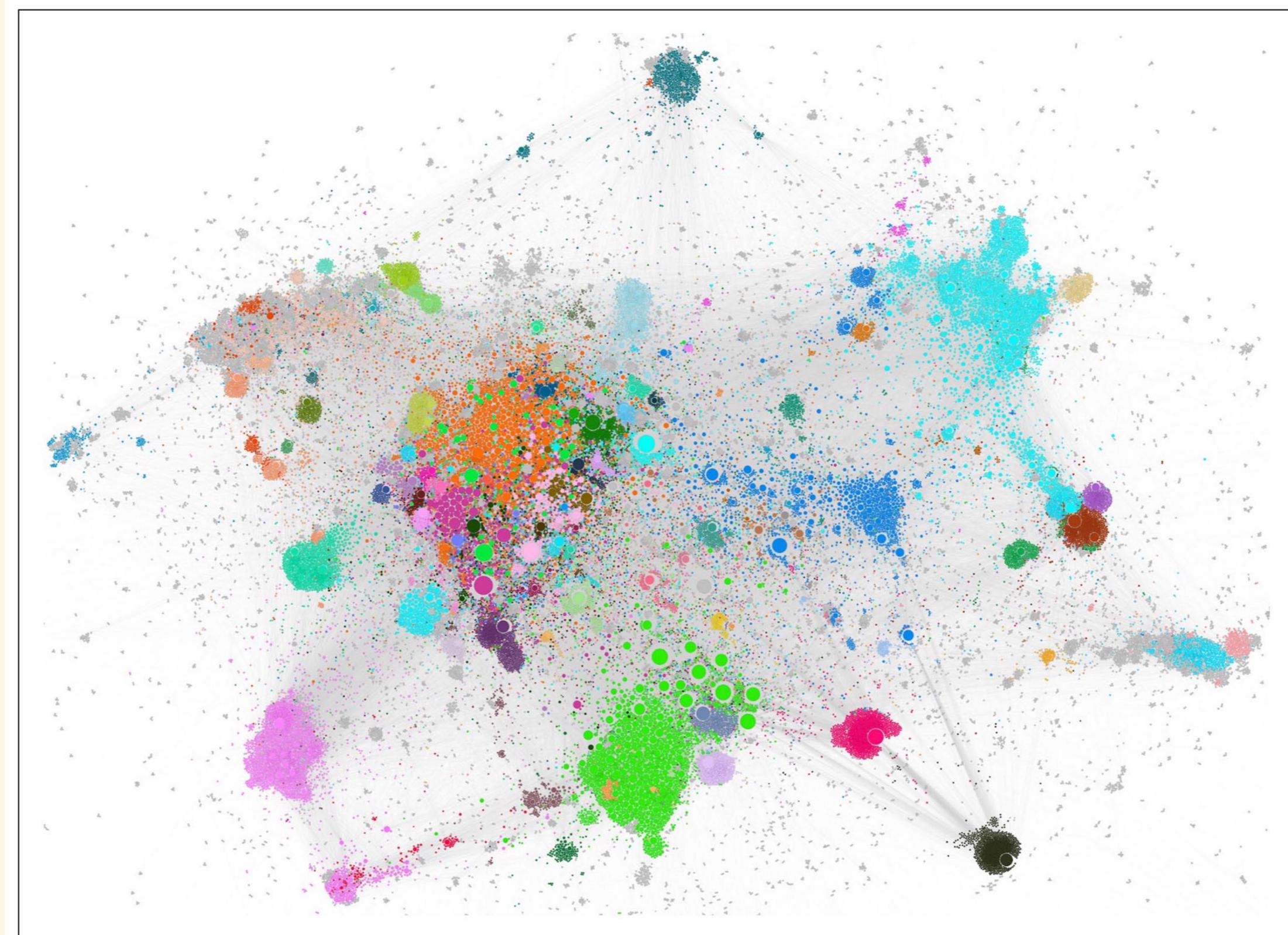


Figure 12. Central communities in the 3-core of our sample network; colored by largest communities detected with the Infomap community detection algorithm (Rosvall & Bergstrom, 2008; Rosvall et al., 2009); node size represents Page Rank (Brin & Page, 1998); layout done with Force Atlas 2 in Gephi (Bastian et al., 2009); (colored version available online).

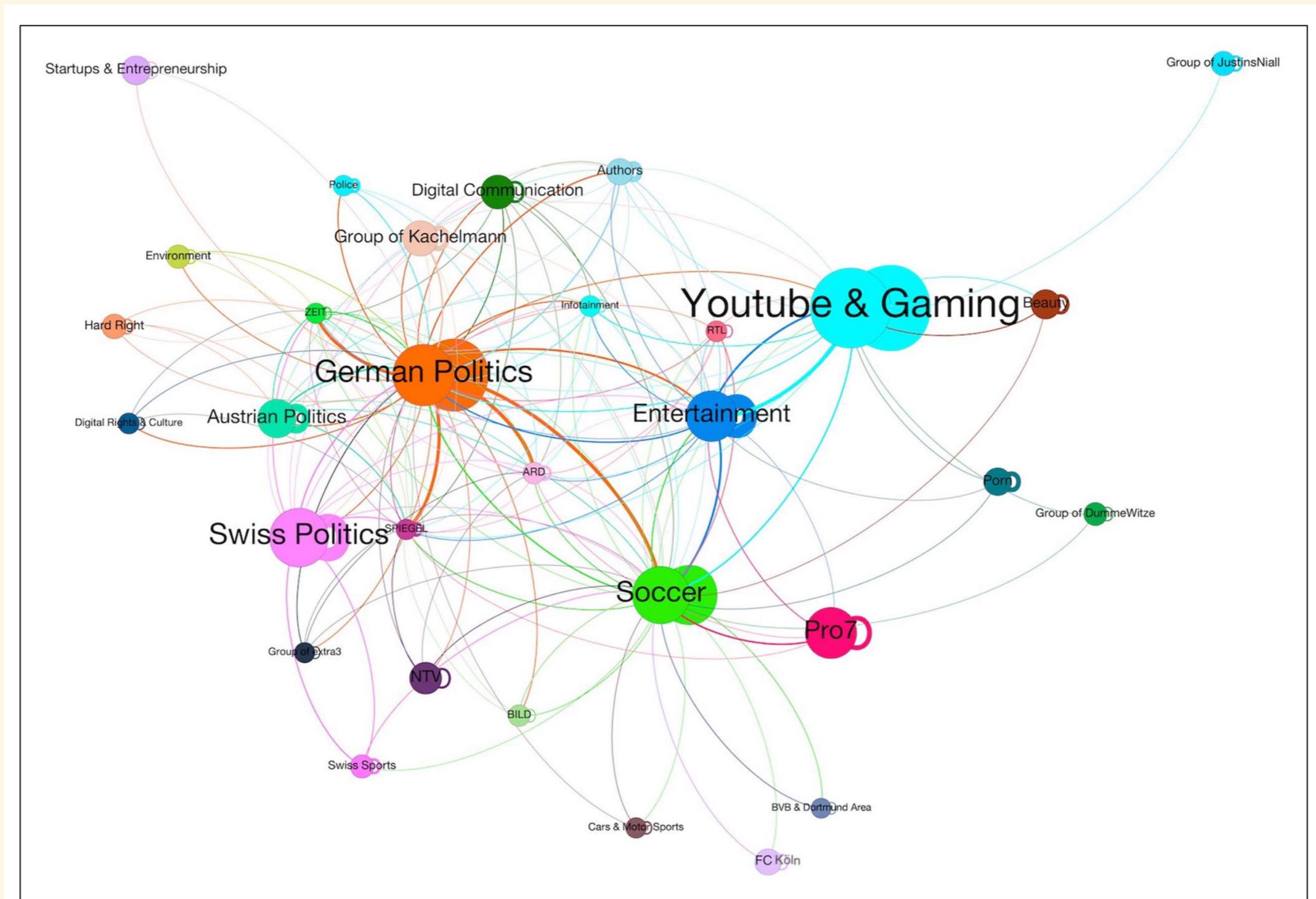


Figure 13. Community graph of communities in the 3-core of our sample with over 300 accounts, at least 80 active accounts during the examined timeframe, and edges with a weight of at least 150; edge width represents weight; edge direction follows clockwise curvature; edges colored by source node; node size represents the number of accounts in each community; node colors correspond with Figure 12; node labels based on interpretation of keywords and top accounts (see Supplemental Material); (colored version available online).

GERMAN-ITALIAN TWITTERSPHERE

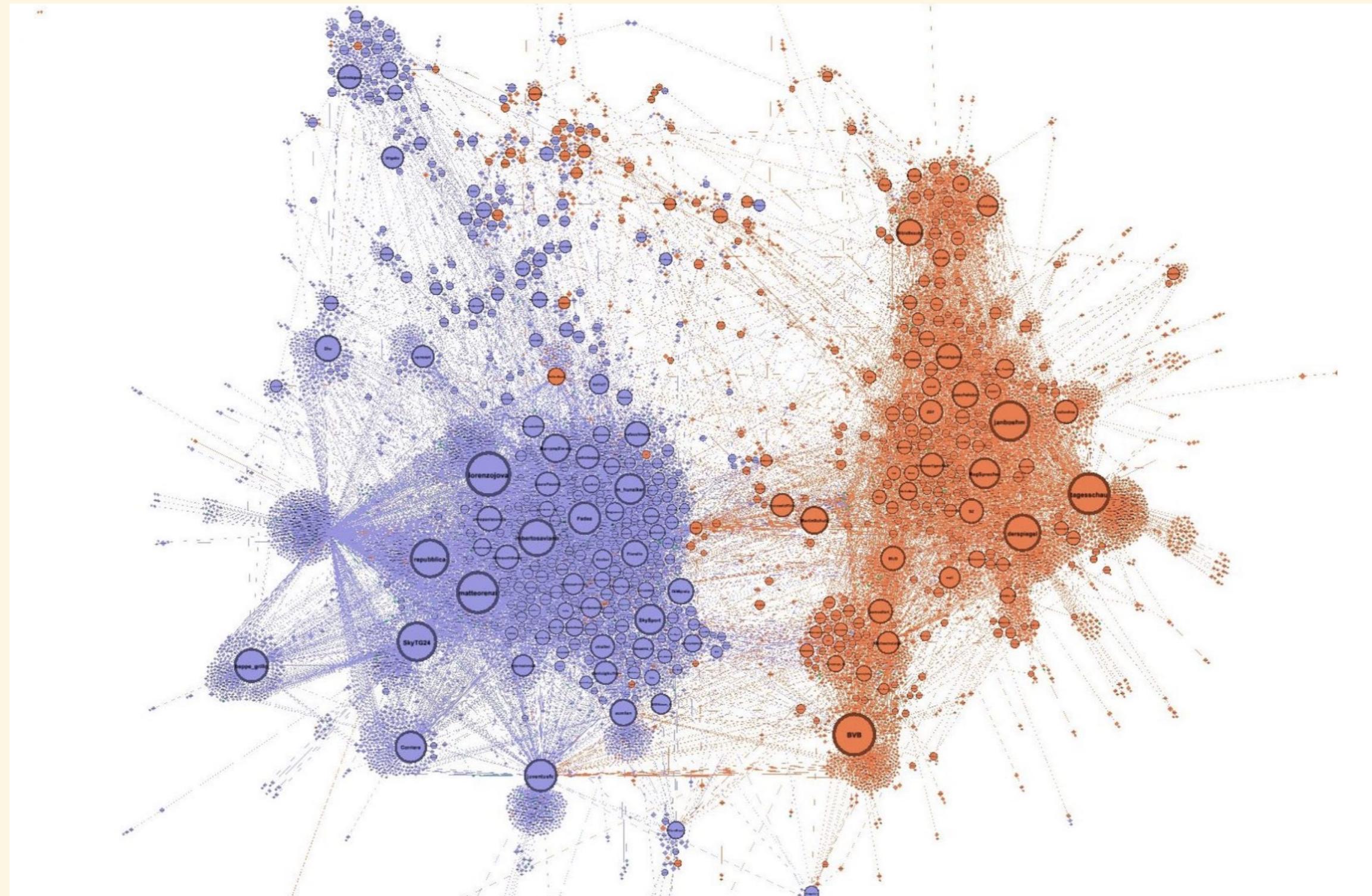


Figure 1: Force directed layout (Force Atlas 2 (Bastian et al., 2009)) of the Italian (purple) - German (orange) follow network sample. Nodes sized by betweenness centrality (Brandes, 2001).

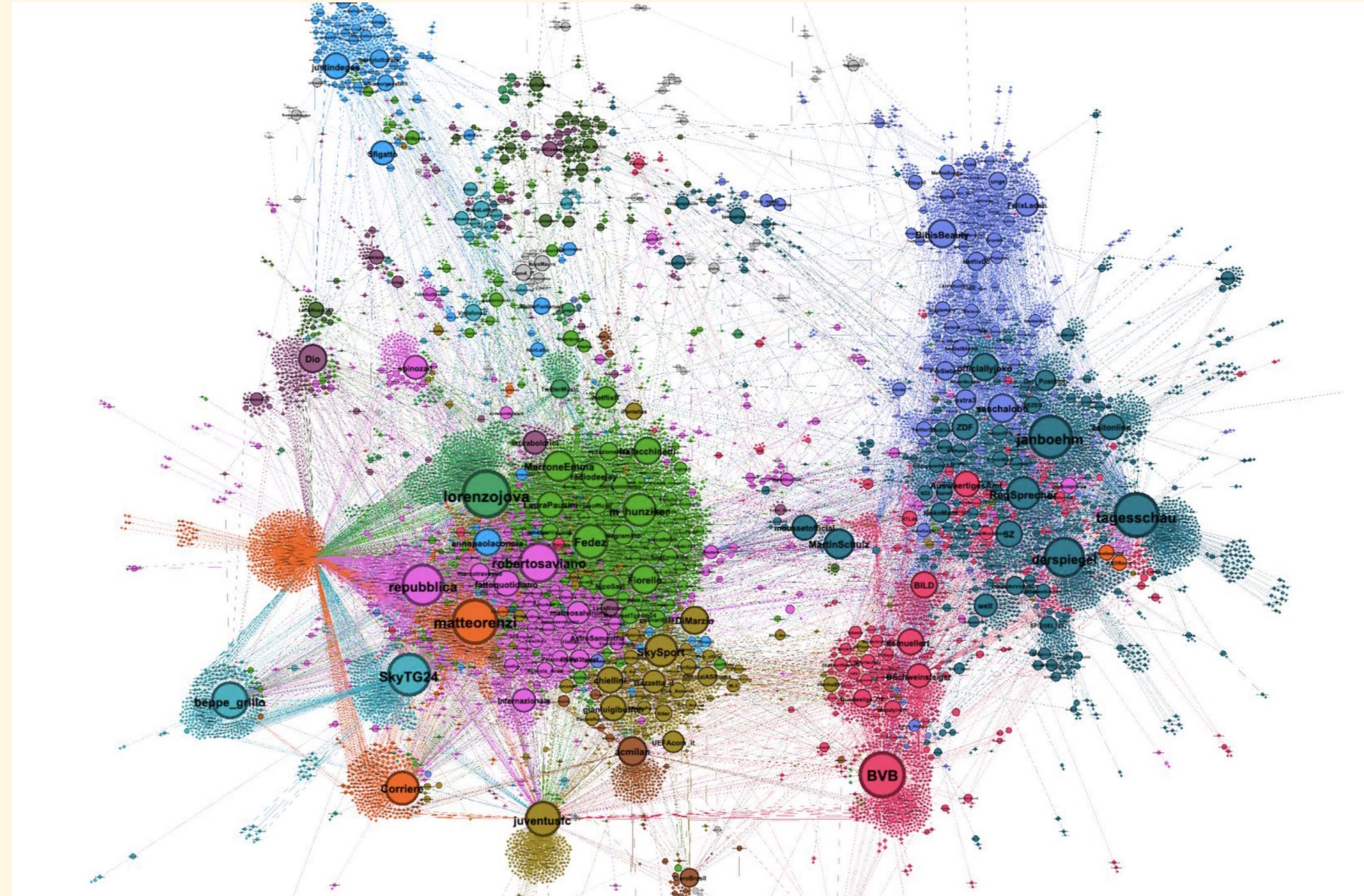


Figure 3: Force directed layout (Force Atlas 2 (Bastian et al., 2009)) of the Italian (left) - German (right) follow network sample. Nodes sized by betweenness centrality (Brandes, 2001). Coloured by modularity maximising clusters.

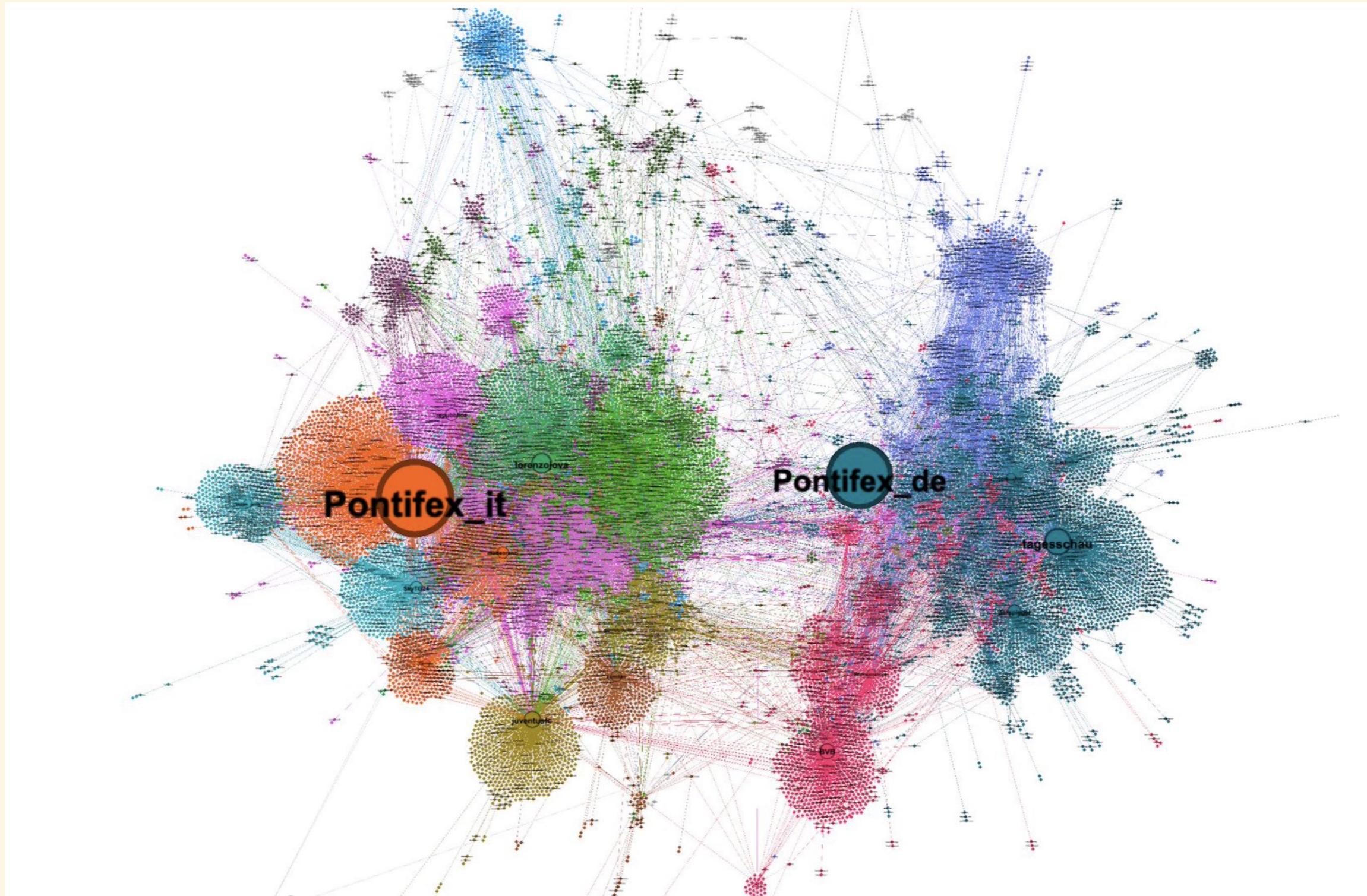
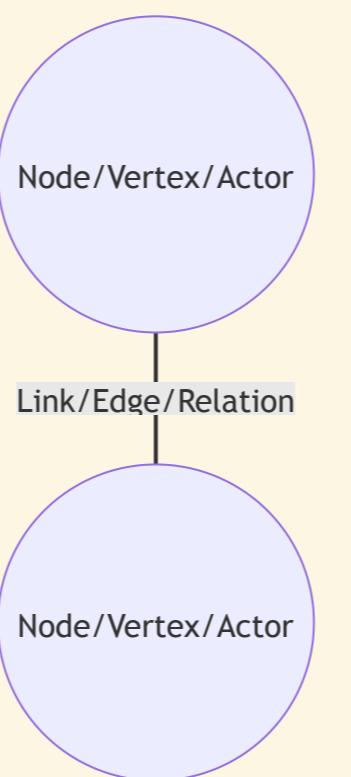


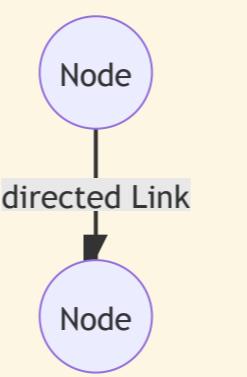
Figure 2: Force directed layout (Force Atlas 2 (Bastian et al., 2009)) of the Italian (left) - German (right) follow network sample. Nodes sized by Page Rank (Brin & Page, 1998). Coloured by modularity maximising clusters.

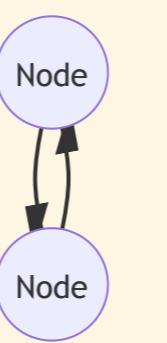
NETWORK ANALYSIS FUNDAMENTALS

ELEMENTS AND PROPERTIES OF NETWORKS

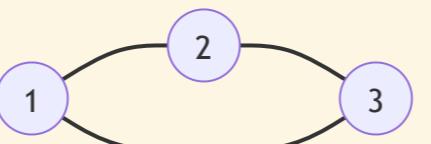
DYADS



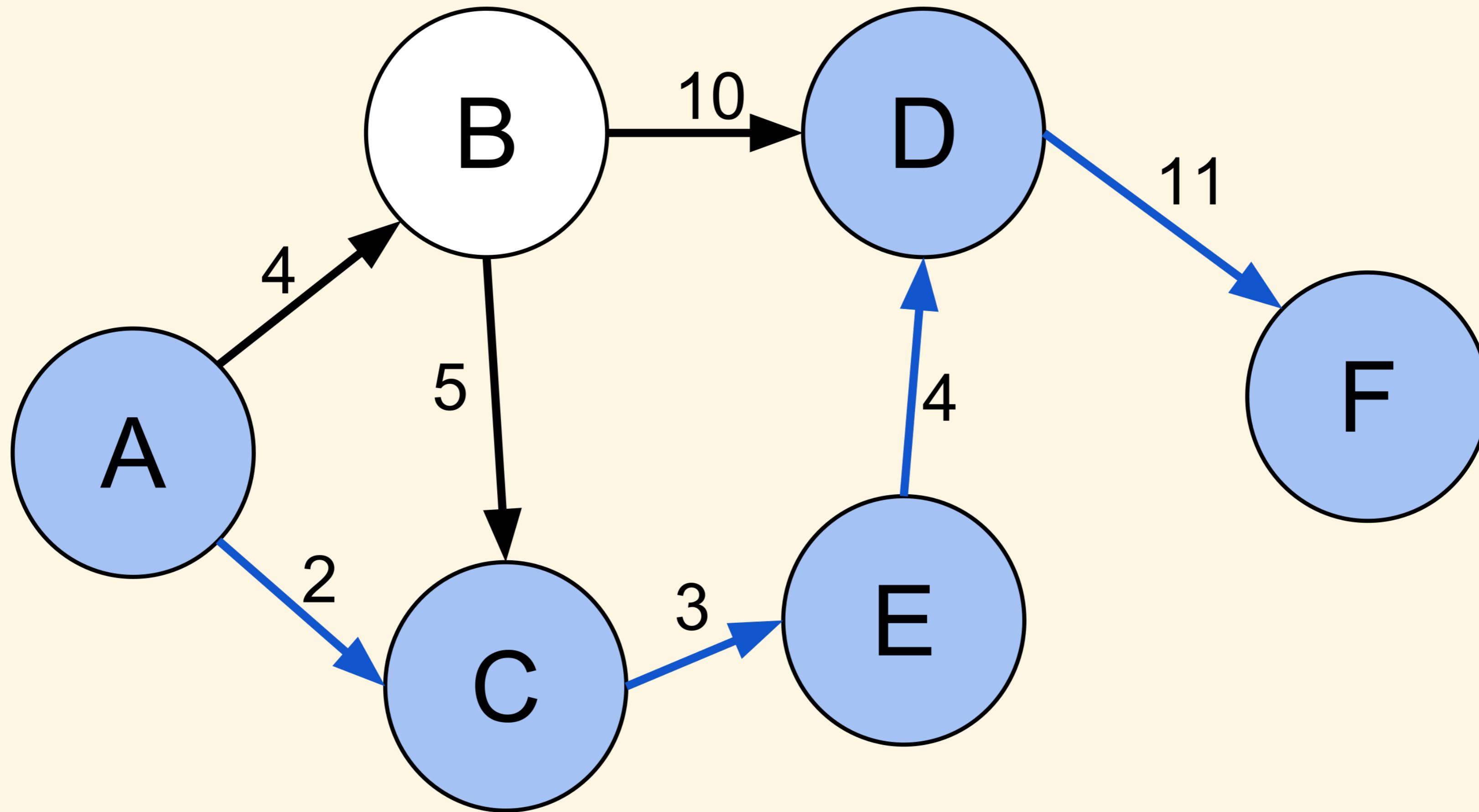




TRIADS



WEIGHTED LINKS & (SHORTEST) PATHS

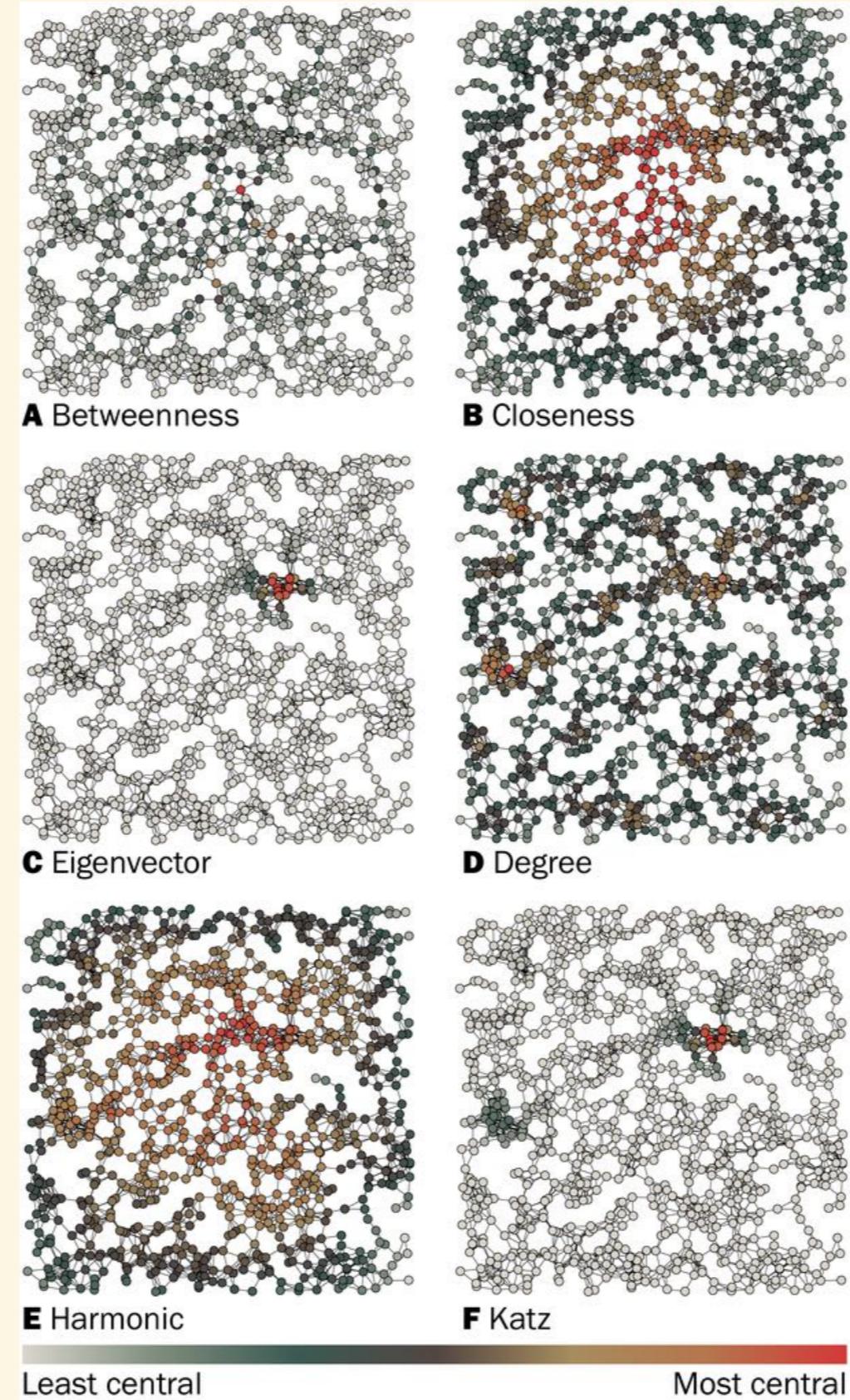


NETWORK ANALYSIS METHODS

MEASUREMENTS OF NETWORKS/GRAPHS AND THEIR ELEMENTS

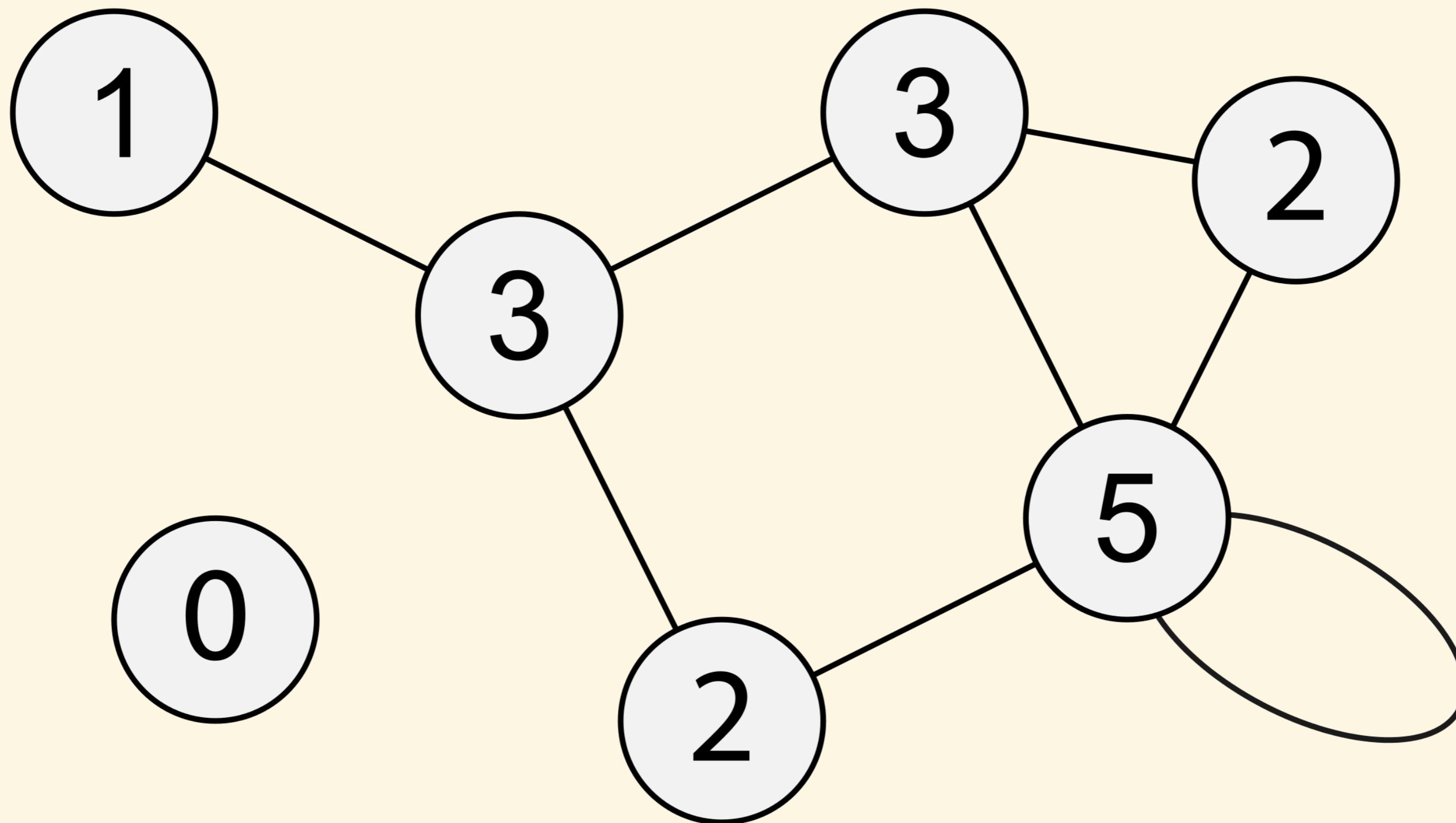
NODE MEASURES

IMPORTANT CENTRALITY MEASURES



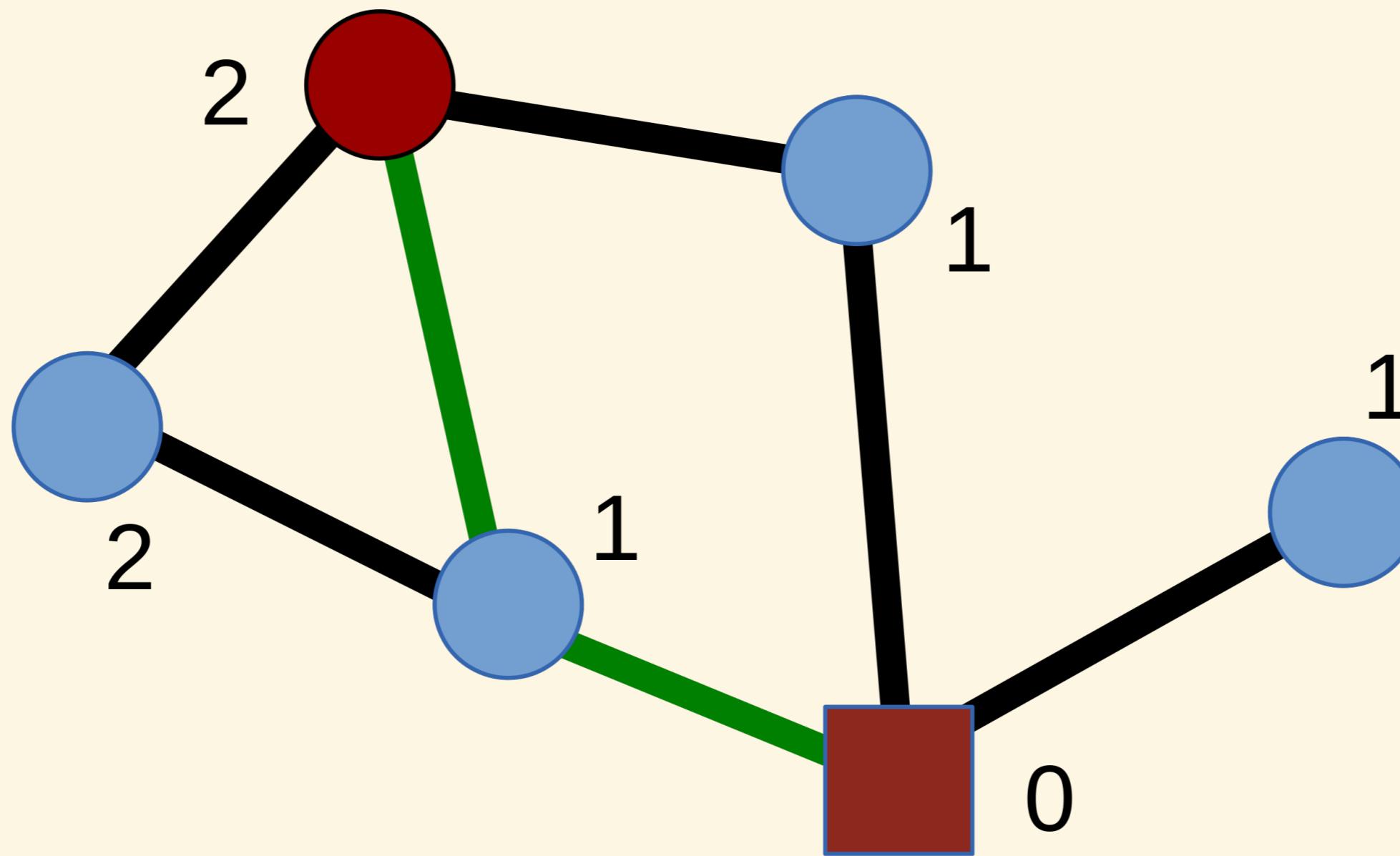
CC-BY-SA Pholme, <https://commons.wikimedia.org/wiki/File:Wp-01.png>

(IN/OUT-)DEGREE CENTRALITY



CLOSENESS CENTRALITY

$$C(x) = \frac{N-1}{\sum_y d(y,x)}$$

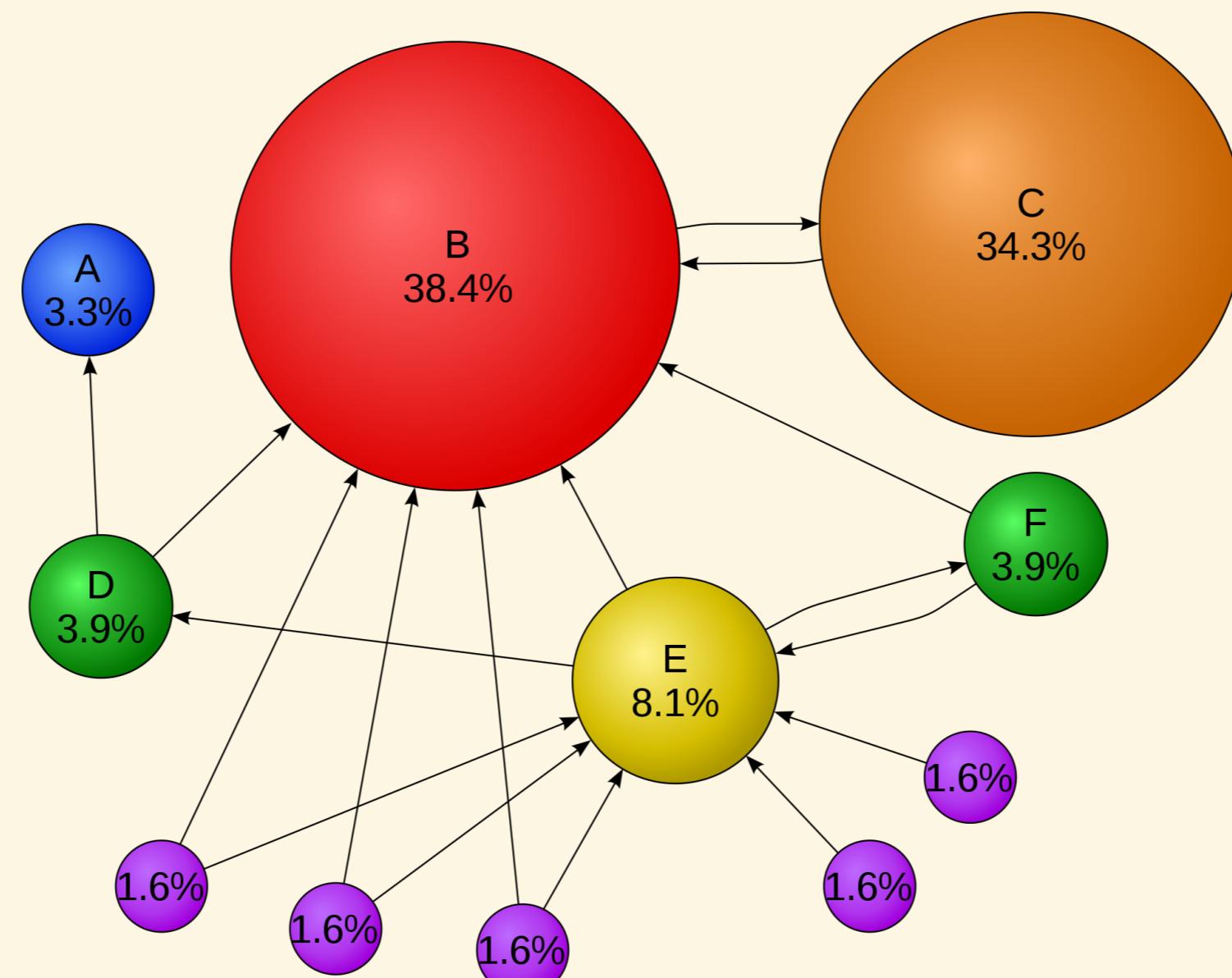


The red square node has (normalised) closeness centrality $\frac{(6-1)}{1+1+1+2+2}$

Image CC-BY-SA CentralConcept, <https://commons.wikimedia.org/wiki/File:Pathdegreeclosenessexampleedit.svg>

EIGENVECTOR CENTRALITY AND PAGE RANK

Both based on the so-called Eigenvalue/Eigenvector equation of the adjacency matrix. In non-math terms:
nodes who have many high ranking nodes as neighbours rank high.

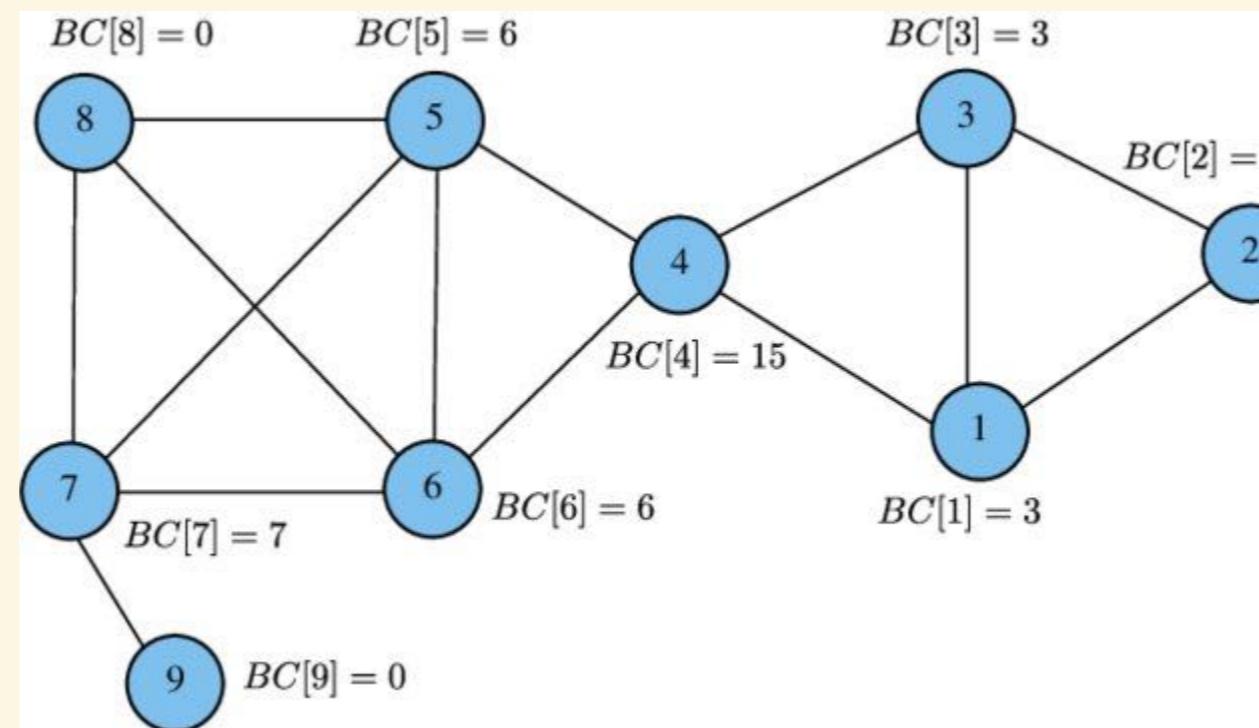


Simple Illustration of the PageRank Algorithm

Image Public Domain, <https://commons.wikimedia.org/wiki/File:PageRanks-Example.svg>

BETWEENNESS-CENTRALITY

$$g(v) = \sum_{s \neq v \neq t} \frac{\sigma_{st}(v)}{\sigma_{st}}$$

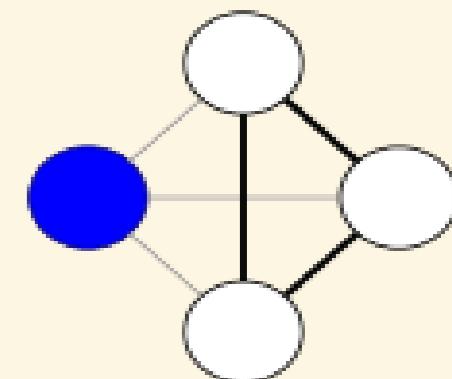


The more shortest paths are going through a node, the higher its betweenness.

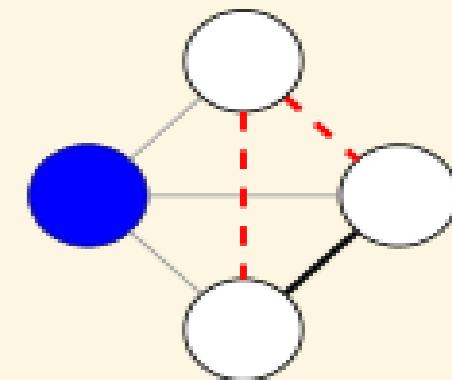
Image Source: McLaughlin, Adam & Bader, David. (2015). Scalable and High Performance Betweenness Centrality on the GPU. International Conference for High Performance Computing, Networking, Storage and Analysis, SC. 2015. 572-583. 10.1109/SC.2014.52.

LOCAL CLUSTERING COEFFICIENT

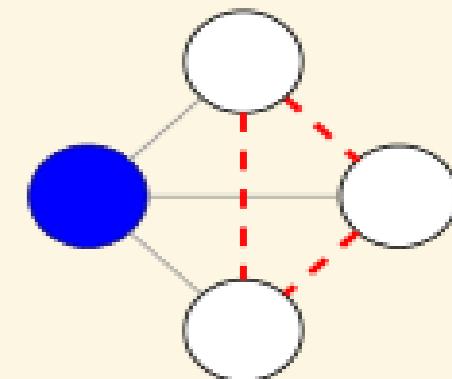
The number of realized edges divided by the number of possible edges between neighbouring nodes



$$c = 1$$



$$c = 1/3$$



$$c = 0$$

IMPORTANT GLOBAL NETWORK MEASURES

Global Clustering Coefficient = $\frac{\text{number of closed triplets}}{\text{number of all triplets}}$

Diameter: longest shortest path of the network

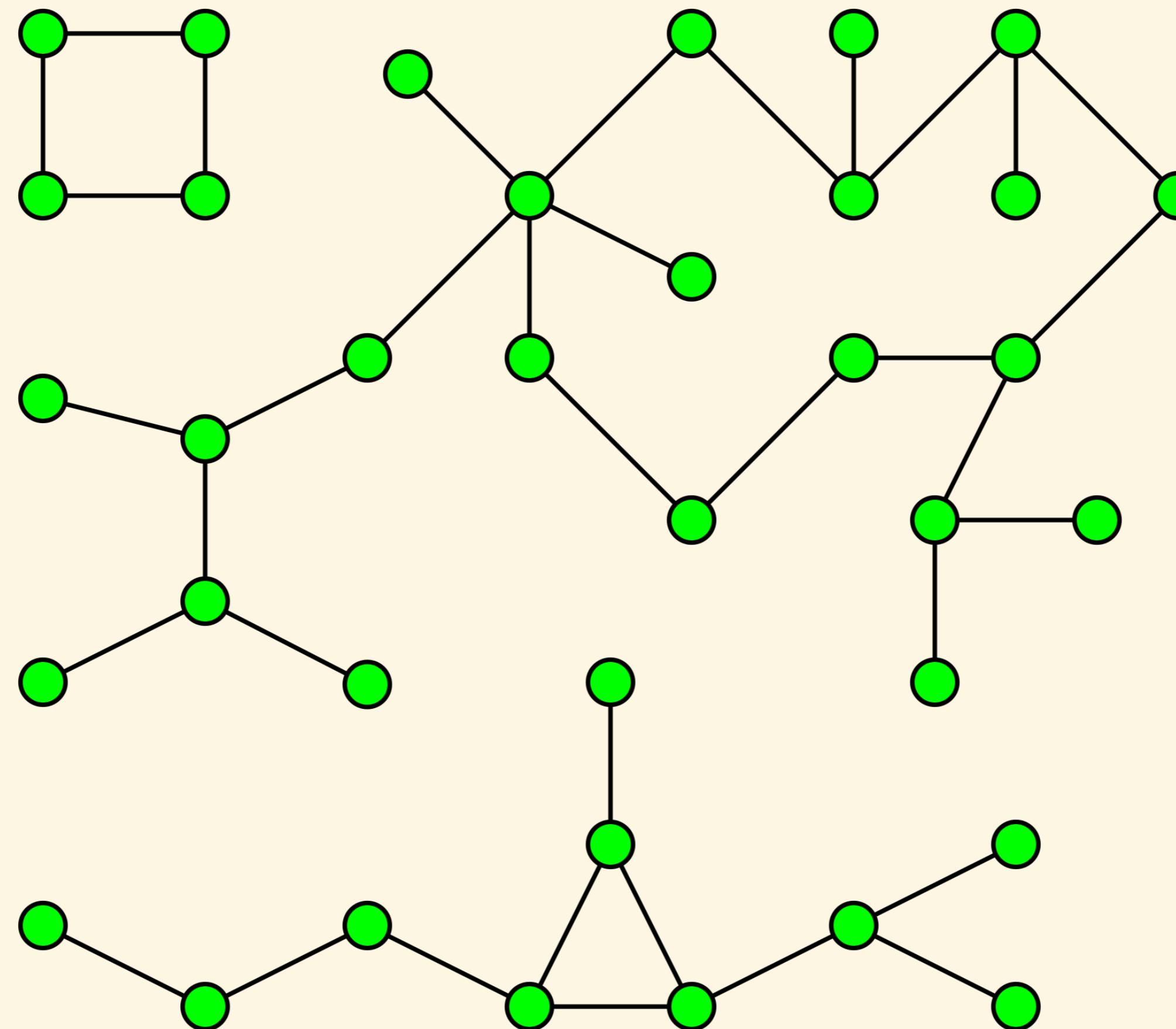
Density: = $\frac{\text{number of links}}{\text{number of possible links}}$

Average Shortest Path Length

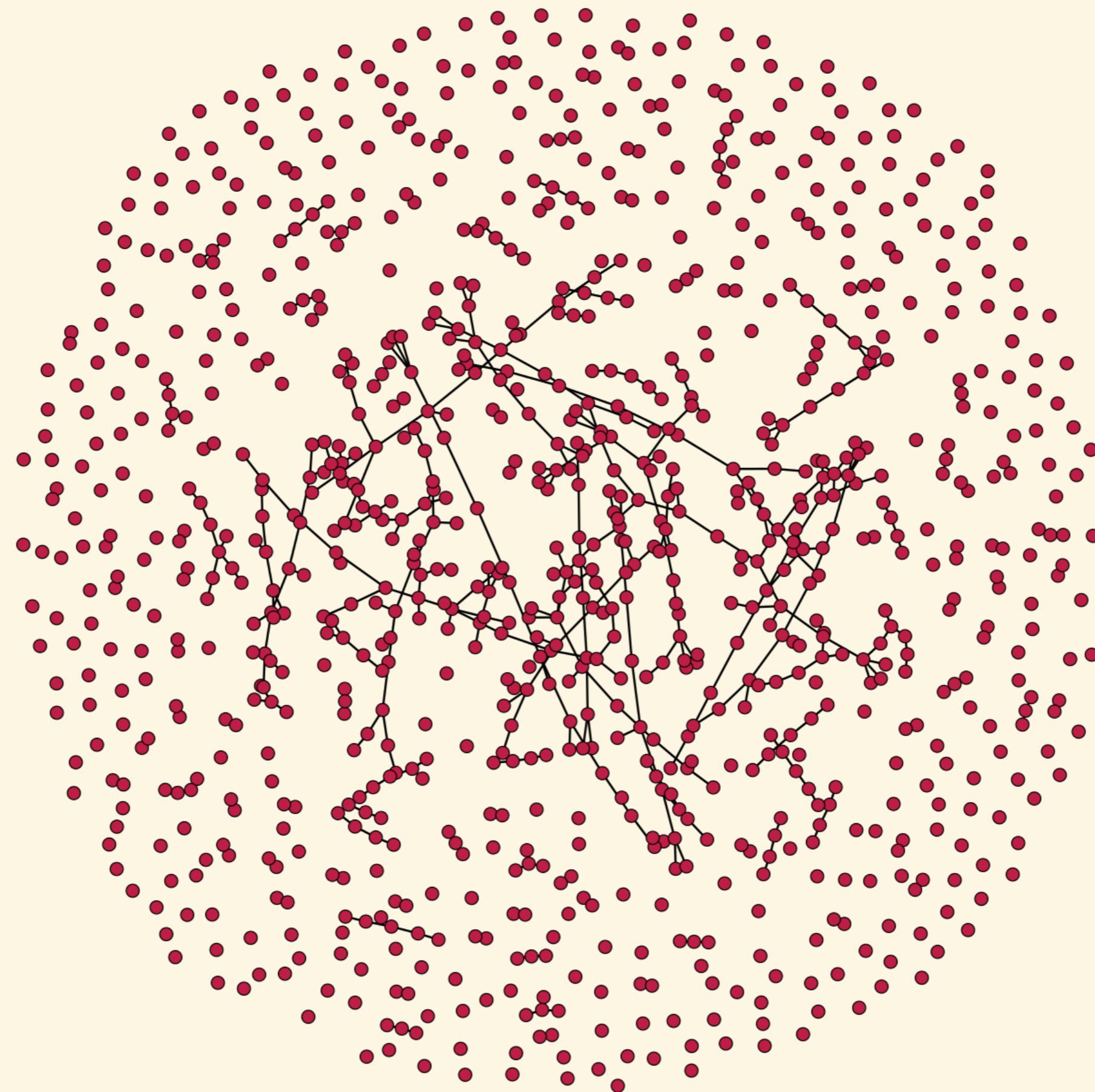
& Averages of most node measures (e.g., average degree, betweenness, closeness, ...)

NETWORKS WITHIN NETWORKS

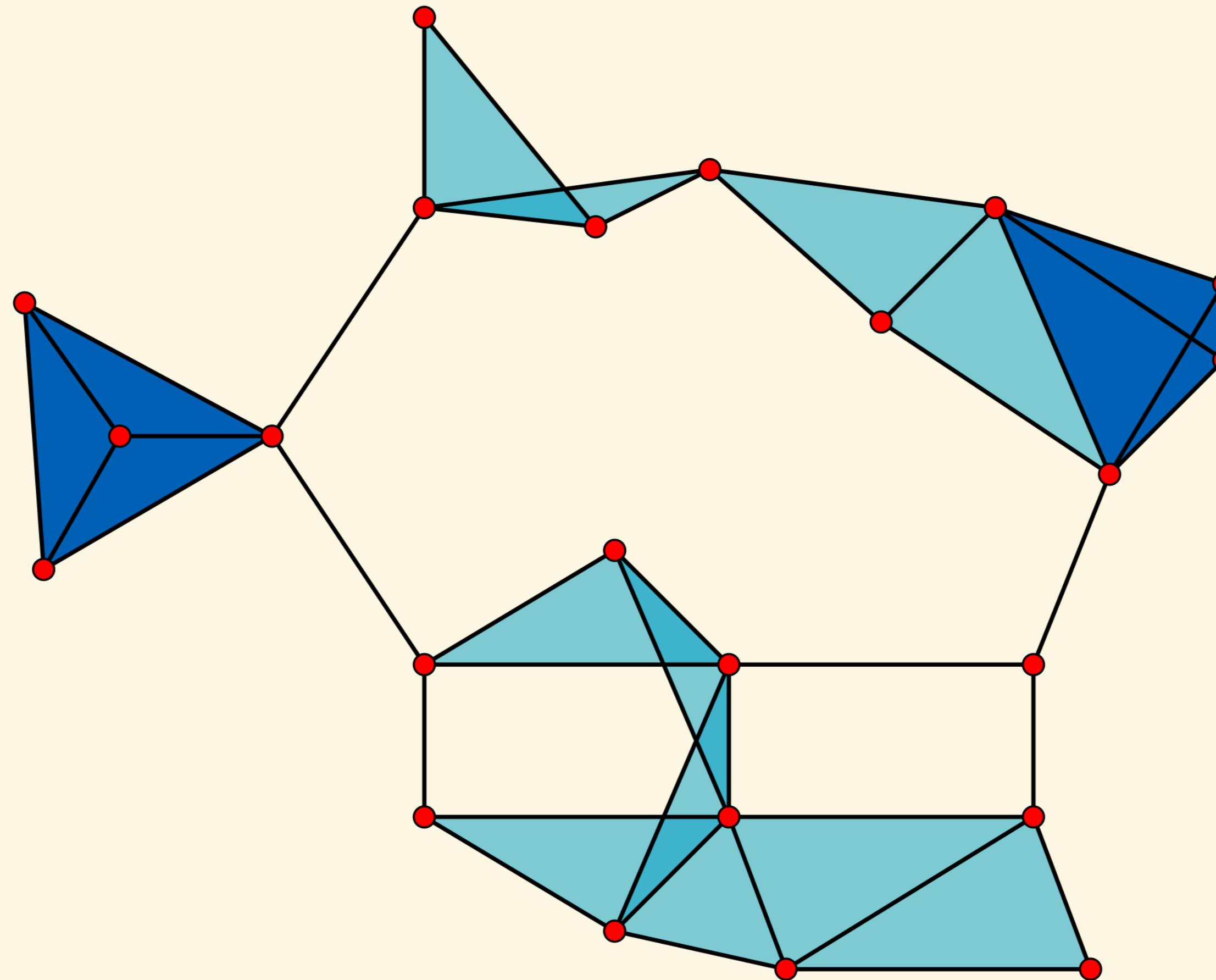
(WEAKLY) CONNECTED COMPONENTS



GIANT COMPONENT

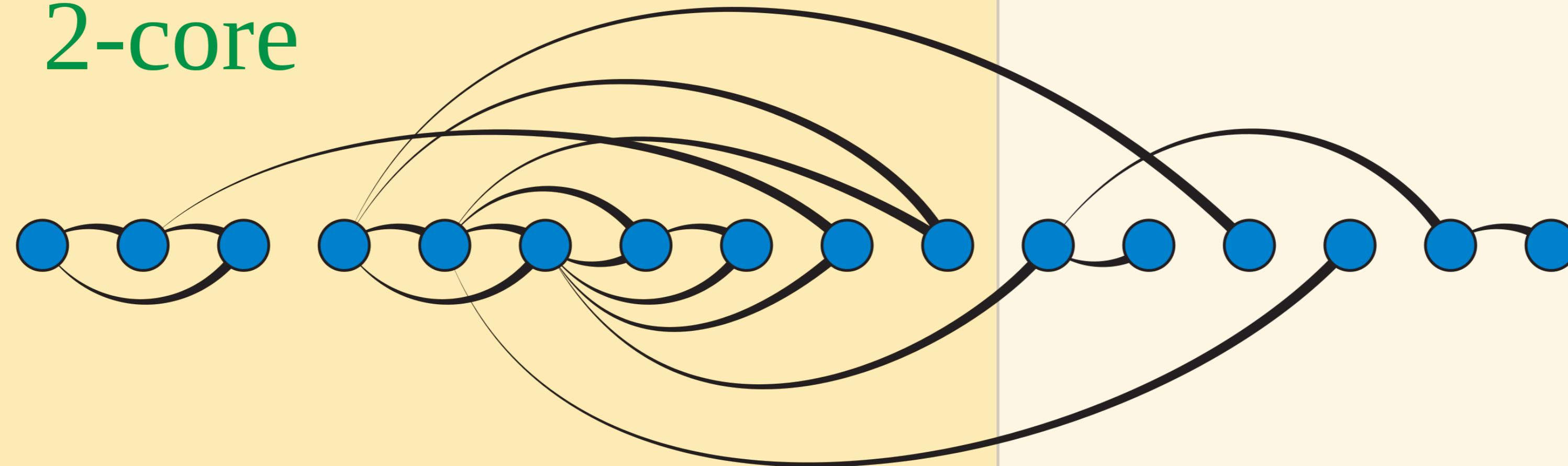


CLIQUEs



K-CORES

2-core



(can also be used as a node centrality measure)

COMMUNITIES/CLUSTERS

Depend on the detection algorithms used. Two of the most popular are

- **Modularity Maximisation (mostly in the Louvain implementation)** based on the relative density of in-/out-group edges
 - and
- **Map Equation (infomap)** based on the length of stay of random walks in certain regions of the network (technically the minimization of the description length of its path)

FLAT COMMUNITIES

- simplify complex systems
- easy to understand
- can oversimplify
- inherently have a resolution limit/arbitrary resolution

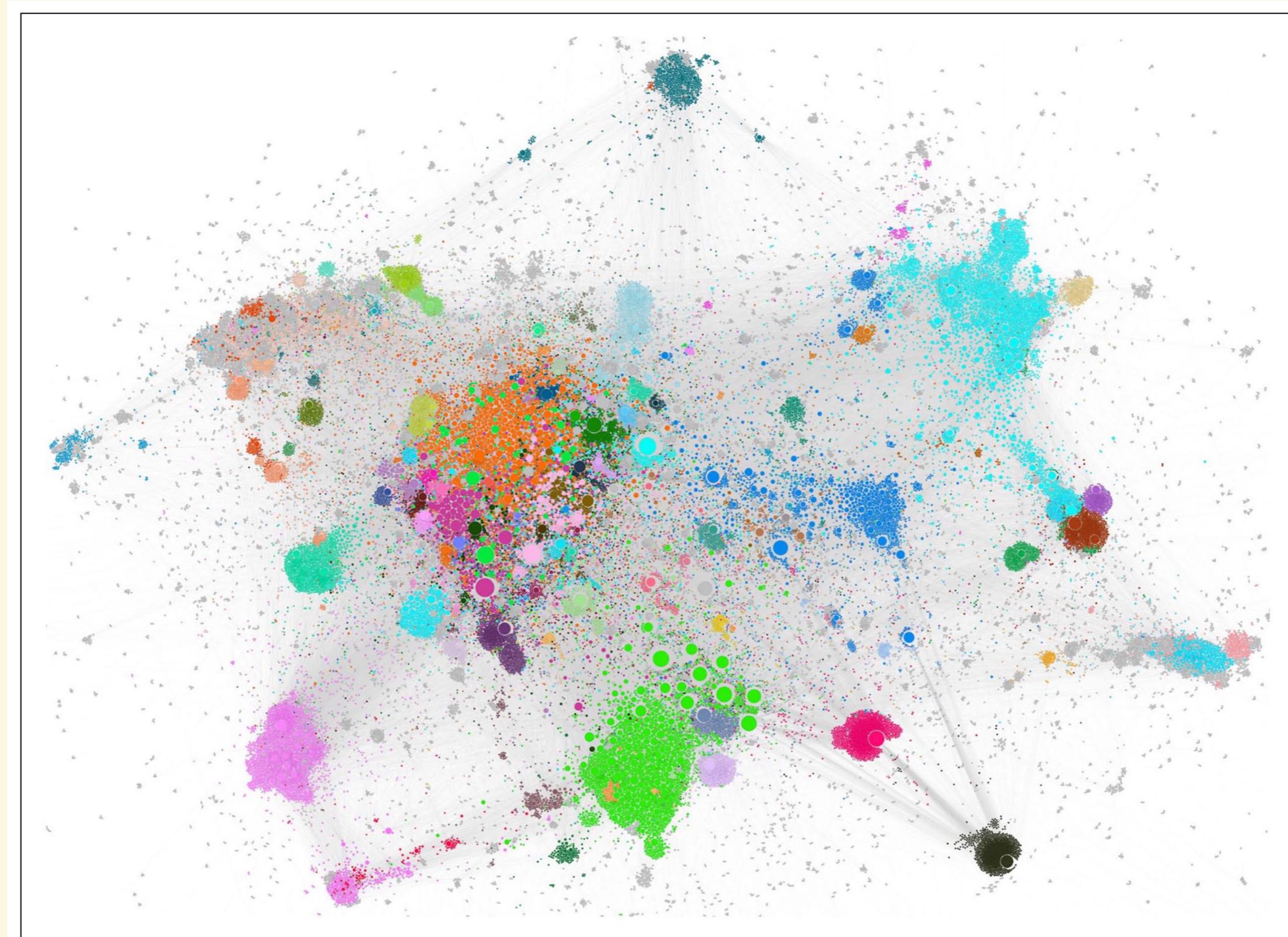
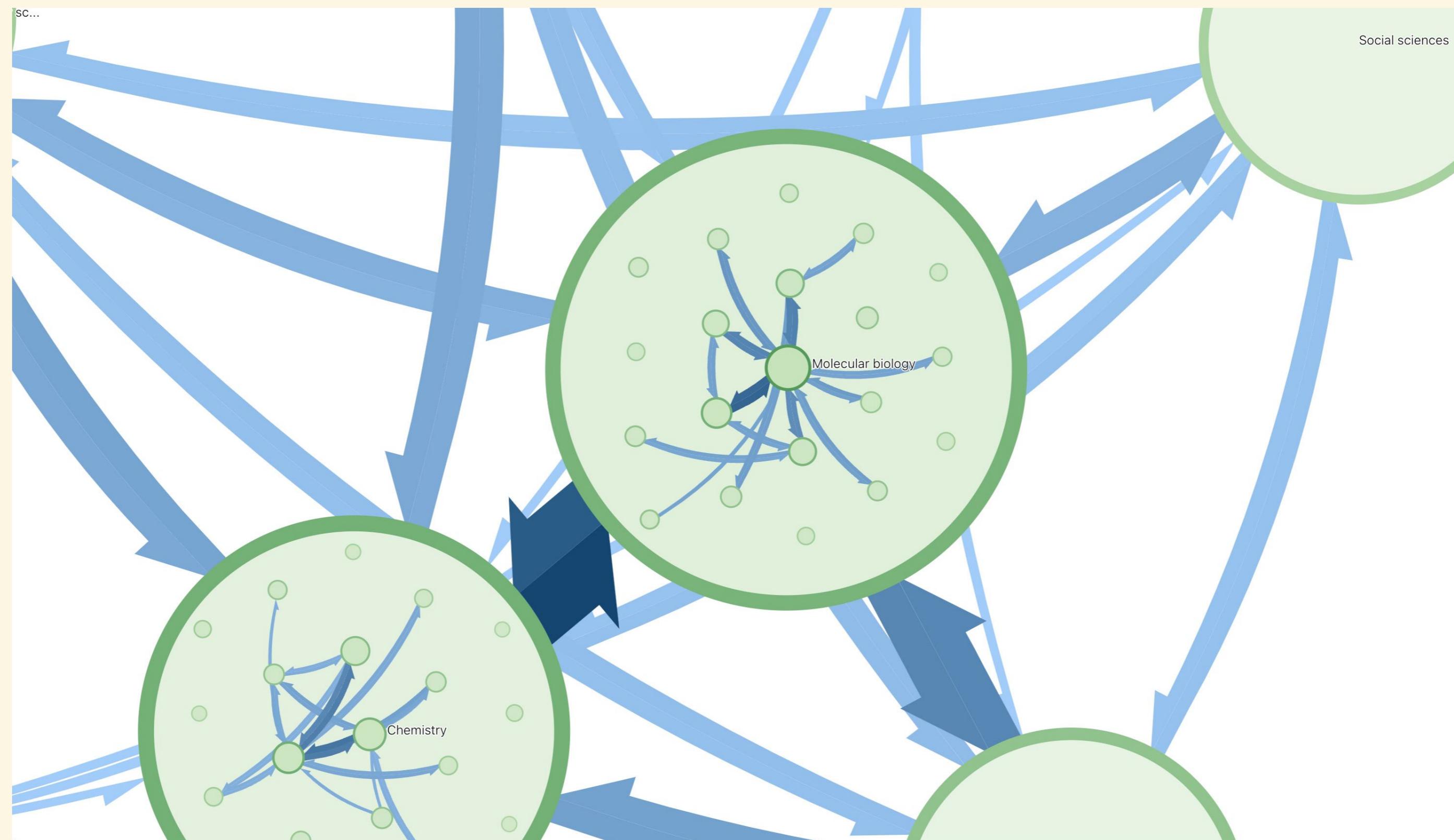


Figure 12. Central communities in the 3-core of our sample network; colored by largest communities detected with the Infomap community detection algorithm (Rosvall & Bergstrom, 2008; Rosvall et al., 2009); node size represents Page Rank (Brin & Page, 1998); layout done with Force Atlas 2 in Gephi (Bastian et al., 2009); (colored version available online).

HIERARCHICAL AND OVERLAPPING COMMUNITIES

- Possible, e.g. with infomap
- often closer to reality
- often easier to interpret
- harder to analyse and visualise



<https://www.mapequation.org/navigator/>

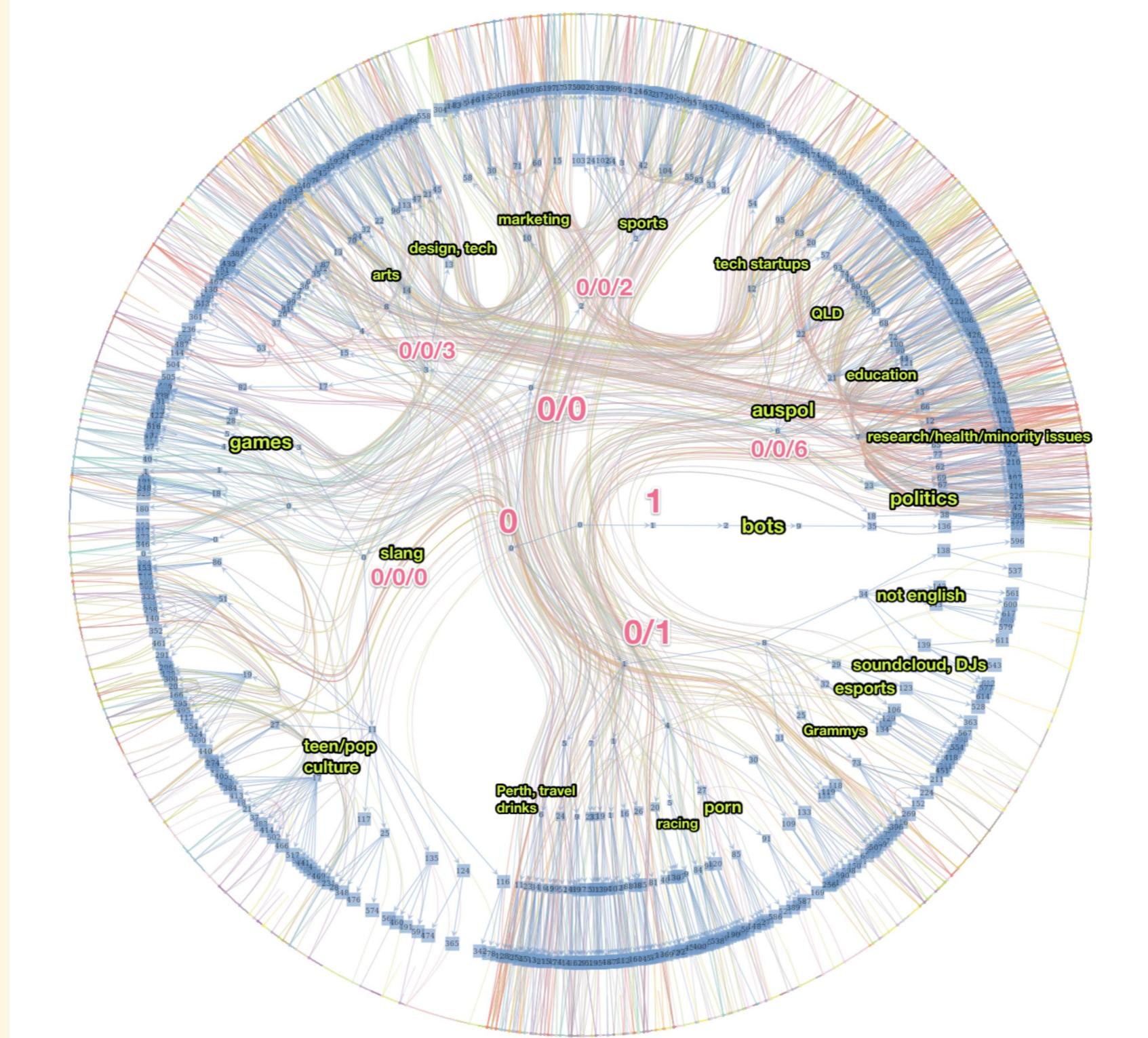


Figure 6.22: Annotated visualisation of the nested stochastic block model inferred for the filtered Australian follow network, edges sampled down to 1000 edges; nodes on outer circle represent accounts; labels based on clearly interpretable results from keyword extraction for blocks detected on level 2; labelling not complete and for overview only; pink labels used for references in text. (High resolution version available online.)

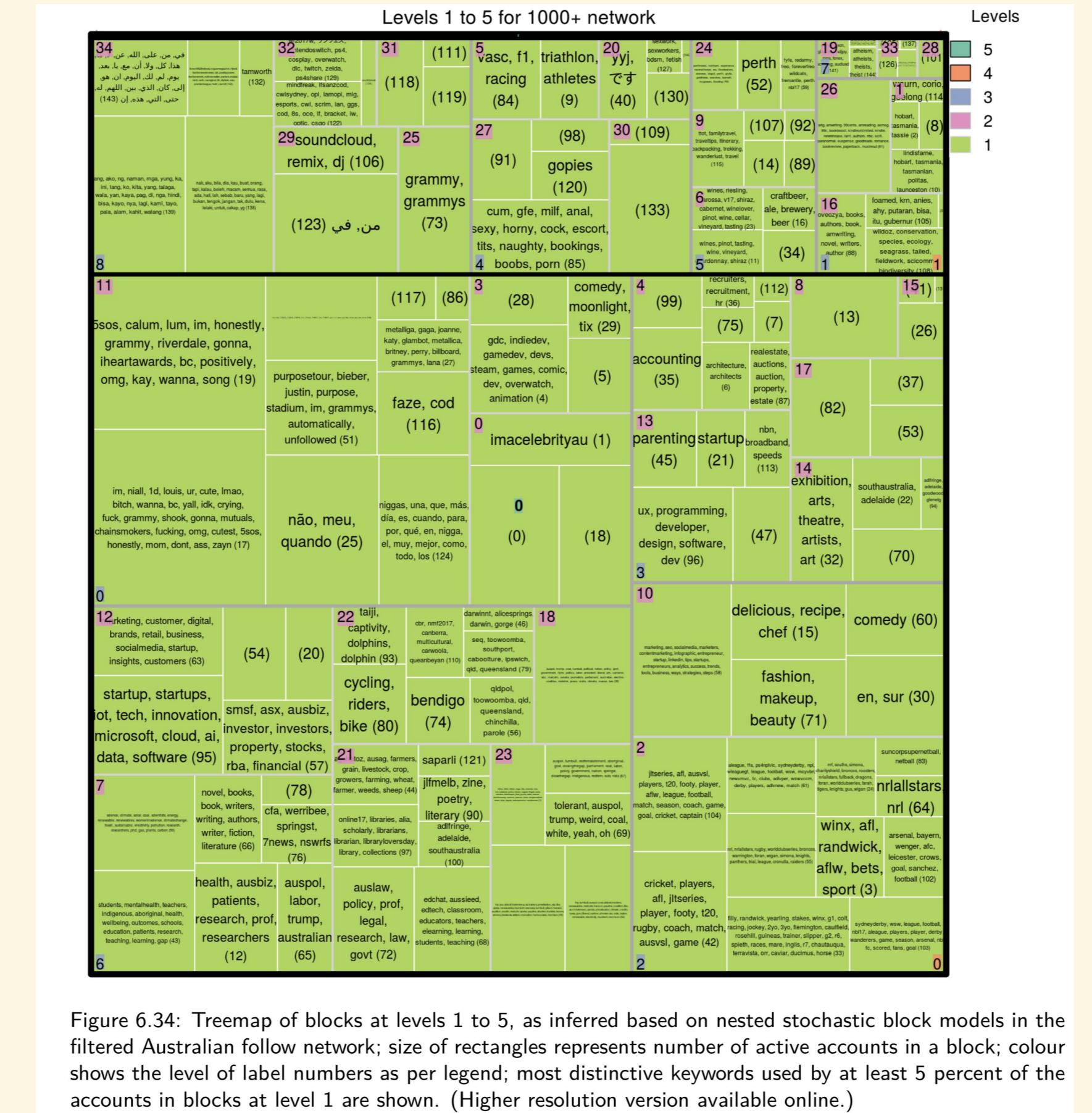


Figure 6.34: Treemap of blocks at levels 1 to 5, as inferred based on nested stochastic block models in the filtered Australian follow network; size of rectangles represents number of active accounts in a block; colour shows the level of label numbers as per legend; most distinctive keywords used by at least 5 percent of the accounts in blocks at level 1 are shown. (Higher resolution version available online.)

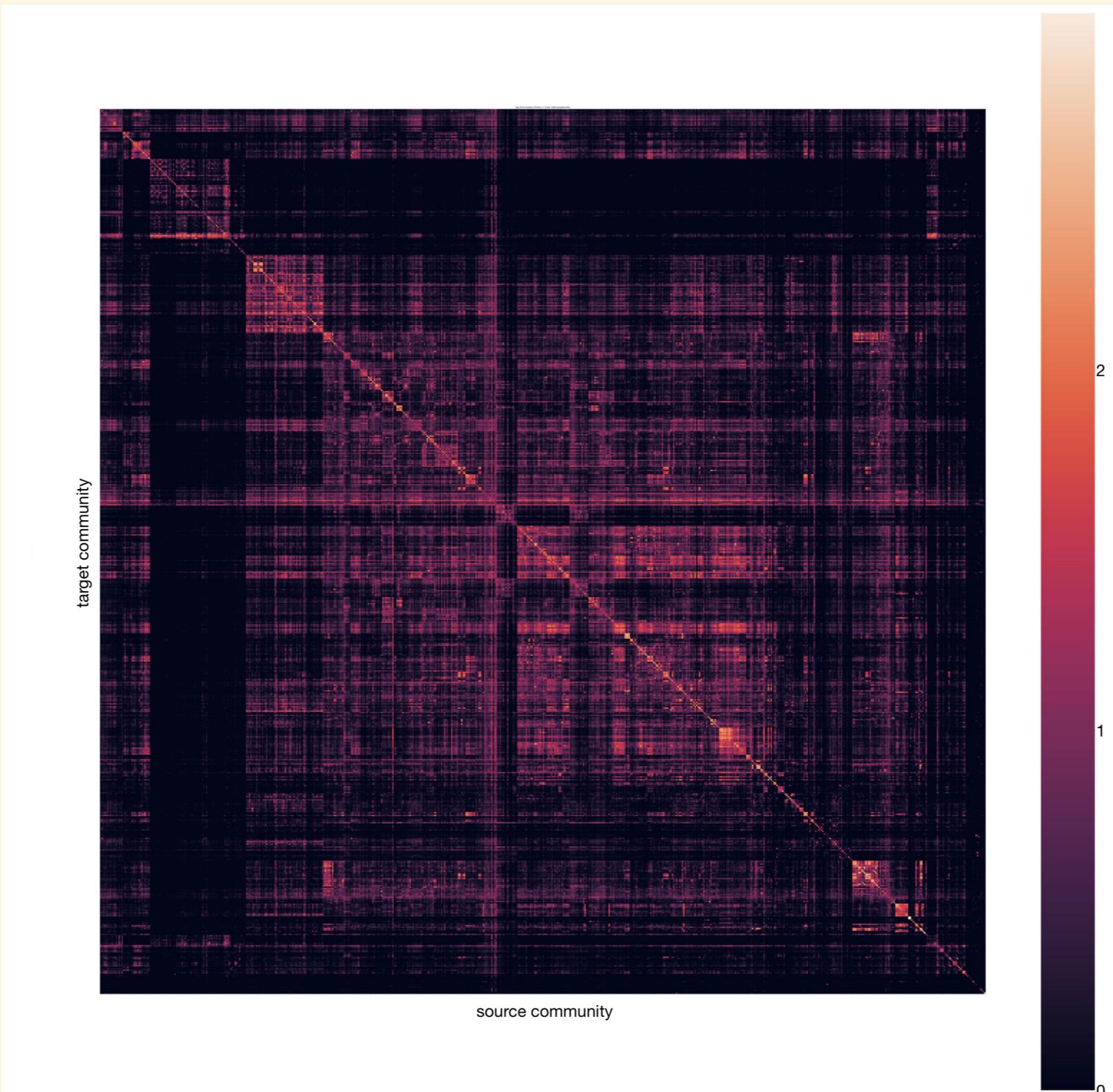


Figure 6.40: Visualisation of the adjacency matrix of the directed block graph of the filtered Australian follow network at level 0, as inferred based on nested stochastic block models; colour shows the logarithm to base 10 of the number of connections + 1 per 1000 possible connections between the source block on the horizontal, and the target block on the vertical axis.

DATA SOURCES FOR (ONLINE ((SOCIAL) MEDIA)) NETWORKS

REPOSITORIES

e.g.:

- Netzschleuder: <https://networks.skewed.de/>,
- SNAP datasets: <https://snap.stanford.edu/data/index.html>
- Network Repository: <https://networkrepository.com/>
- and many more

REPOSITORIES

Pro	Contra
Easy to access	Old data
Fewer legal and ethical problems	Already studied, harder to find new topics
Good for meta studies	Need to trust the data collector
Good for method testing	Available info not tailored to your needs

API

Pro

New/live data

More control over what to collect

Relatively stable machine readability

Legally quite safe

Sometimes access to additional metadata

Contra

Often not historical data

Ethics and data protection considerations apply

Often vetting and acceptance of Terms of Service (TOS)

Rate limits and accessible data shape research question

Can be deprecated/shut down

WEBSCRAPING

Pro

New/live data

More control over what to collect

No vetting or acceptance of TOS

No rate limits

In best case WYSIWYG

Contra

Kind of unstable machine readability

More Ethics and data protection considerations apply

Active countermeasures by platforms

Technically often more complex setup

Also possible, but less common for large networks:

- Surveys
- Questionnaires
- Data Donations
- "Manual" Collection
- ...

QUESTIONS?

@flxvctr

