

DECODING SOCIAL INTERACTIONS IN THE WEB

[Rossano Schifanella]

@eBISS 2018

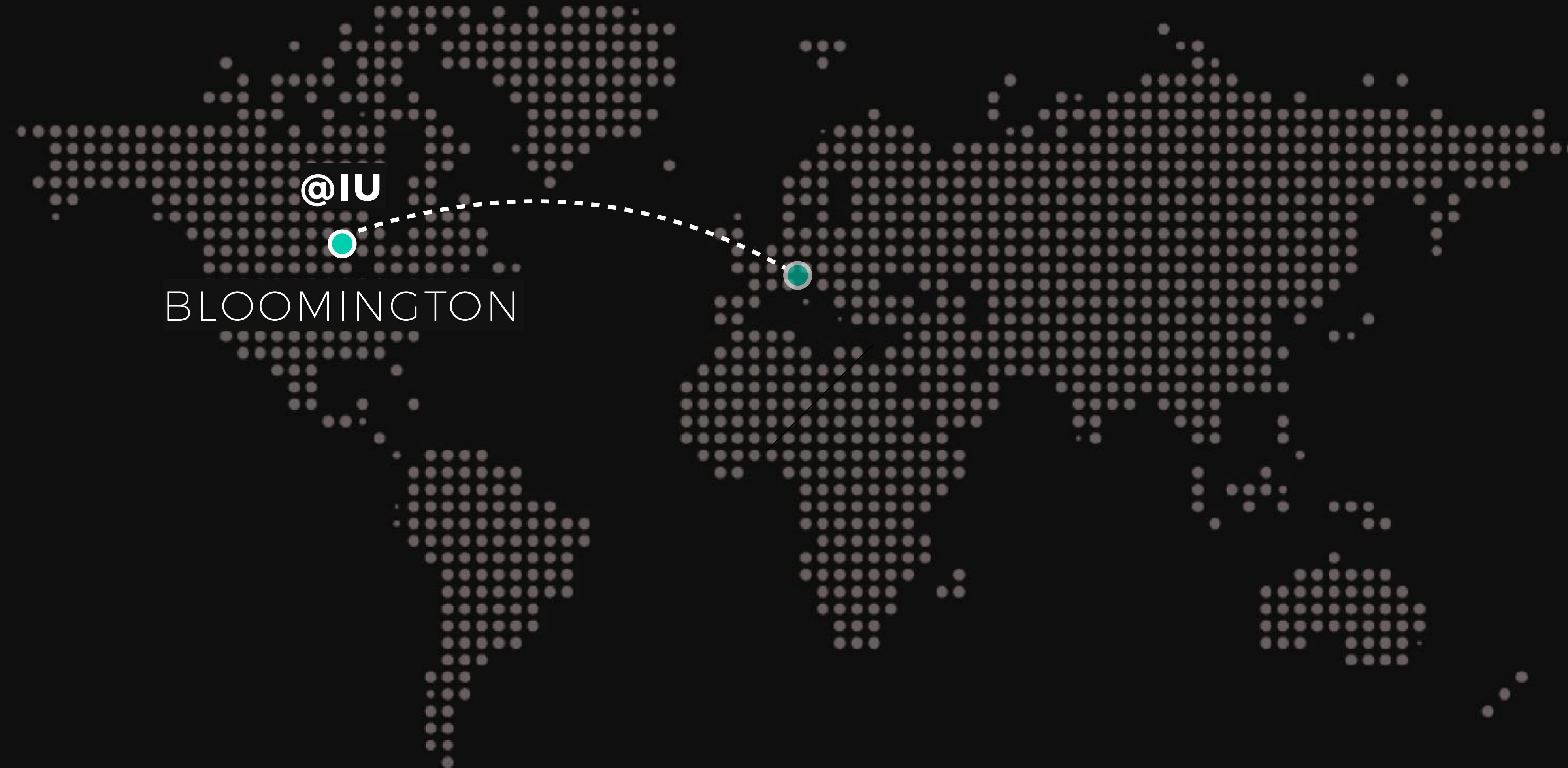
Berg en Dal, The Netherlands

Who Am I?



@UNITO
●
TURIN

DISTRIBUTED SYSTEMS
RECOMMENDER SYSTEMS



NETWORK SCIENCE
HUMAN COMPUTING

CROWDSOURCING



URBAN COMPUTING
MACHINE LEARNING

BEHAVIOURAL STUDIES
COMPUTATIONAL *



@Nokia Bell Labs

CAMBRIDGE

URBAN COMPUTING
COMPUTATIONAL SOCIAL SCIENCE

The image features a world map where each country is represented by a grid of small gray dots. Specific locations are highlighted with teal-colored circles. The labeled locations and their approximate coordinates are:

- BLOOMINGTON (USA) - located in the midwest, marked with a teal circle.
- @IU (Indiana University) - located near Bloomington, with the handle @IU above it.
- NEW YORK (USA) - located on the east coast, marked with a teal circle.
- @Yahoo (Yahoo) - located in New York City, with the handle @Yahoo above it.
- CAMBRIDGE (UK) - located in the northeast, marked with a teal circle.
- @Nokia Bell Labs (Nokia Bell Labs) - located in Cambridge, with the handle @Nokia Bell Labs above it.
- BARCELONA (Spain) - located in the northeast, marked with a teal circle.
- TURIN (Italy) - located in the northwest, marked with a teal circle.
- @UNITO (Università degli Studi di Torino) - located in Turin, with the handle @UNITO above it.

Who Are You?

Acknowledgements

Slides and material have been inspired by:

CSS Tutorial@WWW16

C. Wagner, L. Aiello, M. Strohmaier, Ingmar Weber

Computational Social Science course

Claudia Wagner, Universität Koblenz-Landau

Algorithmic Bias Tutorial@KDD16

H. Hajian, F. Bonchi, C. Castillo

Social Network Course

Dennis M. Feehan, Berkeley (very extensive literature review)

Computational Social Science Course

John McLevey, University of Waterloo

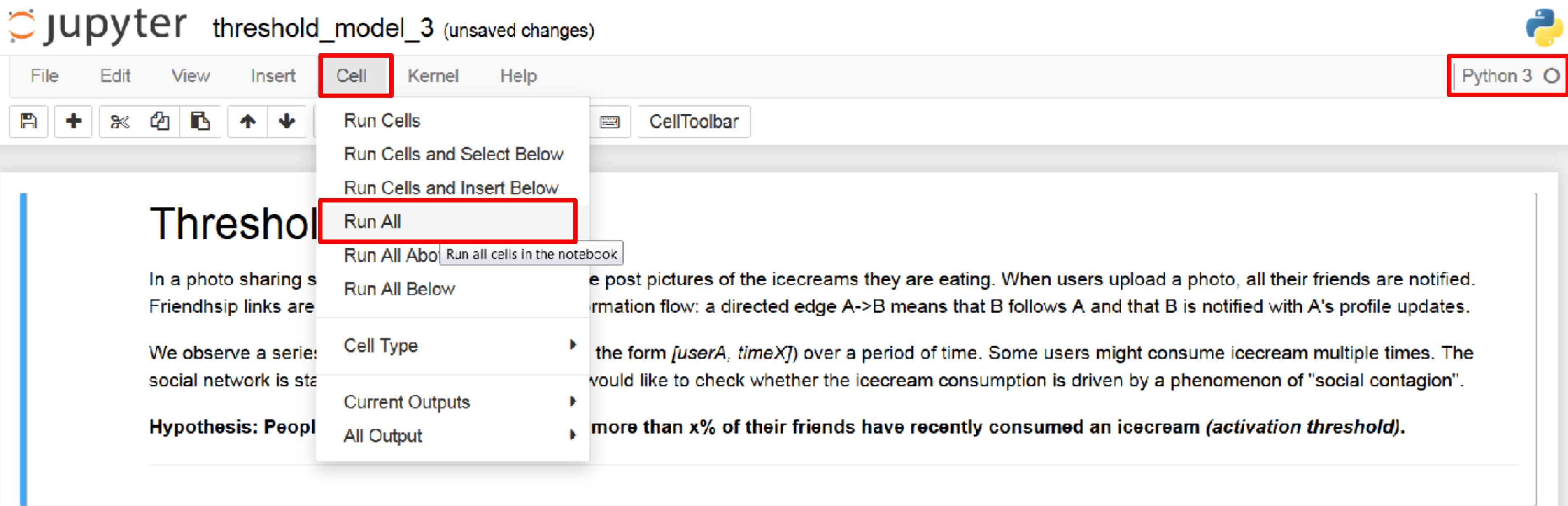
Agenda

- Introduction to Computational Social Science
- The Web as a data source
- Modeling social interactions through networks
- Basic networks metrics
- The importance of social theories

Setup for Practical Session

Setup for Practical Exercises

- We will have practical examples throughout the tutorial
 - Demoing some code. Follow along or try later.
- I have already prepared IPython notebooks for you
- Requirements:
 - Python 3
 - Project Jupyter (<https://jupyter.org/>)
 - networkx
- Download data and source files from:
 - <https://github.com/rschifan/ebiss2018>



STEP 0: Setup

Imports and utility methods. We set our notion of "recent" to 6 hours (as exercise, you can play with this constant)

```
In [1]: #Imports and variables
import sys
import math
import bisect
import heapq
from random import shuffle

import numpy as np
```

Introduction to Computational Social Science

What is Computational Social Science?

Science that investigates **social phenomena** through the medium of **computing** and **algorithmic data processing**.

Go beyond prediction by adding **interpretability** and **explainability**.

Interdisciplinary data science.

biogeographic patterns. Their study, too, is centered on a large database, but in this case it is entirely of living organisms, the marine bivalves. Over 28,000 records of bivalve genera and subgenera from 322 locations around the world have now been compiled by these authors, giving a global record of some 854 genera and subgenera and 5132 species. No fossils are included in the database, but because bivalves have a good fossil record, it is possible to estimate accurately the age of origin of almost all extant genera. It is then possible to plot a backward survivorship curve (8) for each of the 27 global bivalve provinces (9).

On the basis of these curves, Krug *et al.* find that origination rates of marine bivalves in-

creased significantly almost everywhere immediately after the K-Pg mass extinction event. The highest K-Pg origination rates all occurred in tropical and warm-temperate regions. A distinct pulse of bivalve diversification in the early Cenozoic was concentrated mainly in tropical and subtropical regions (see the figure).

The steepest part of the global backward survivorship curve for bivalves lies between 65 and 50 million years ago, pointing to a major biodiversification event in the Paleogene (65 to 23 million years ago) that is perhaps not yet captured in Alroy *et al.*'s database (5, 7). The jury is still out on what may have caused this event. But we should not lose sight of the fact that the steep rise to prominence of many mod-

ern floral and faunal groups in the Cenozoic may bear no simple relationship to climate or any other type of environmental change (10, 11).

References

1. G. G. Metzlerbach *et al.*, *Ecol. Lett.* **10**, 315 (2007).
2. A. Z. Krug, D. Jiborski, J. W. Valentine, *Science* **323**, 753 (2009).
3. P. W. Signor, *Annu. Rev. Ecol. Syst.* **21**, 539 (1990).
4. R. K. Bambach, *Glob. Biogeochem. Cycles* **13**, L1 (1999).
5. J. Alroy *et al.*, *Proc. Natl. Acad. Sci. U.S.A.* **98**, 6261 (2001).
6. A. M. Bush *et al.*, *Paleobiology* **30**, 666 (2004).
7. J. Alroy *et al.*, *Science* **321**, 91 (2008).
8. M. Foote, in *Evolutionary Patterns*, J. D. C. Jackson *et al.*, Eds. (Univ. of Chicago Press, Chicago, IL, 2001), vol. 245, pp. 245–295.
9. M. D. Spalding *et al.*, *Bioscience* **57**, 570 (2007).
10. S. M. Stanley, *Paleobiology* **33**, 1 (2007).
11. M. J. Benton, S. C. Emerson, *Paleobiology* **50**, 23 (2007).

10.1126/science.1169430

SOCIAL SCIENCE

Computational Social Science

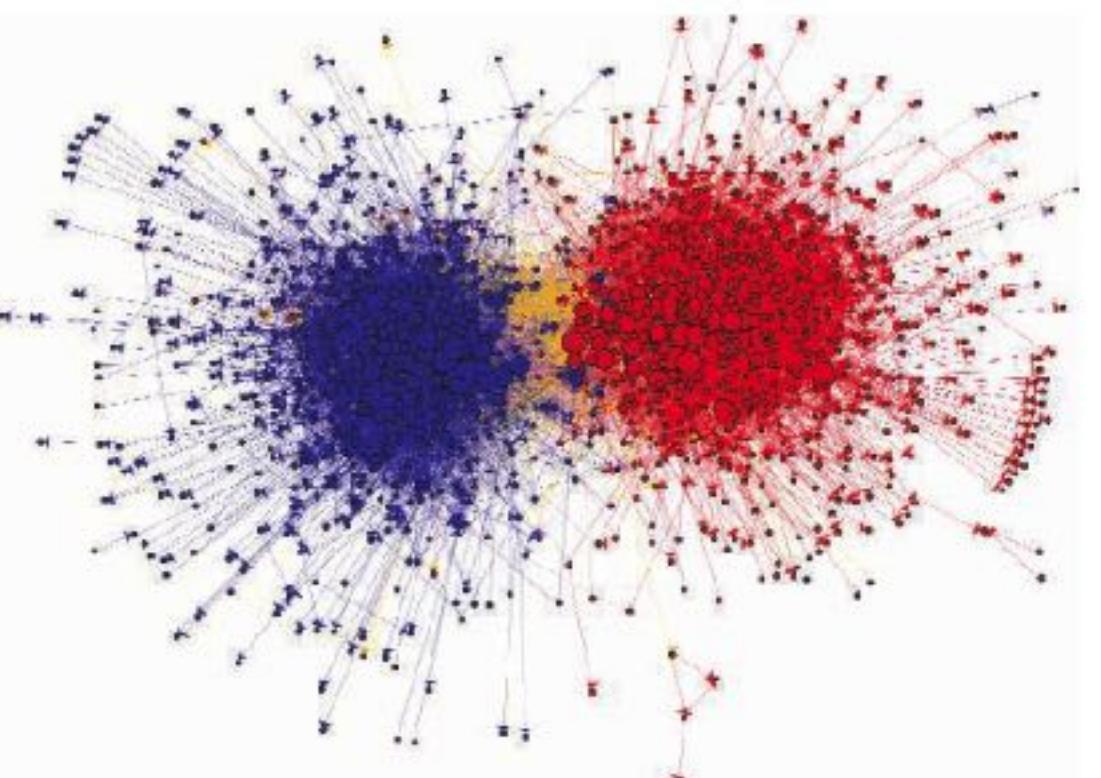
David Lazer,¹ Alex Pentland,² Lada Adamic,³ Sinan Aral,^{2,4} Albert-László Barabási,⁵ Devon Brewer,⁶ Nicholas Christakis,¹ Noshir Contractor,⁷ James Fowler,⁸ Myron Gutmann,³ Tony Jebara,⁹ Gary King,¹ Michael Macy,¹⁰ Deb Roy,² Marshall Van Alstyne^{1,11}

We live life in the network. We check our e-mails regularly, make mobile phone calls from almost any location, swipe transit cards to use public transportation, and make purchases with credit cards. Our movements in public places may be captured by video cameras, and our medical records stored as digital files. We may post blog entries accessible to anyone, or maintain friendships through online social networks. Each of these transactions leaves digital traces that can be compiled into comprehensive pictures of both individual and group behavior, with the potential to transform our understanding of our lives, organizations, and societies.

The capacity to collect and analyze massive amounts of data has transformed such fields as biology and physics. But the emergence of a data-driven “computational social science” has been much slower. Leading journals in economics, sociology, and political science show little evidence of this field. But computational social science is occurring—in Internet companies such as Google and Yahoo, and in govern-

ment agencies such as the U.S. National Security Agency. Computational social science could become the exclusive domain of private companies and government agencies. Alternatively, there might emerge a privileged set of academic researchers presiding over private data from which they produce papers that cannot be critiqued or replicated. Neither scenario will serve the long-term public interest of accumulating, verifying, and disseminating knowledge.

What might a computational social science based in an open academic environment—offer society, by enhancing understanding of individuals and collectives? What are the



Data from the blogosphere. Shown is a link structure with n = a community of political blogs (from 2004), where red nodes indicate conservative blogs, and blue liberal. Orange links go from liberal to conservative, and purple ones from conservative to liberal. The size of each blog reflects the number of other blogs that link to it. [Reproduced from (8) with permission from the Association for Computing Machinery.]

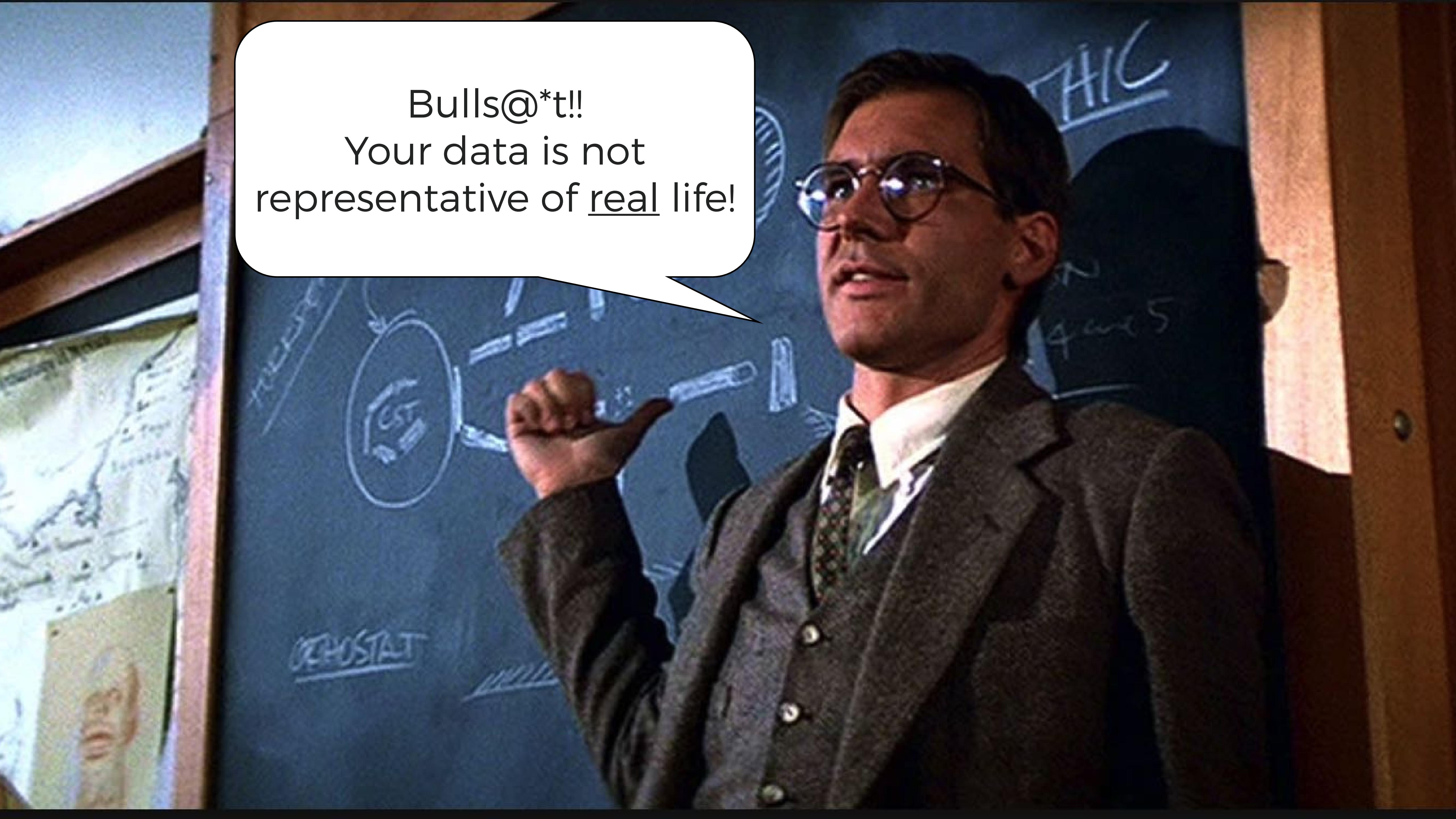


multidisciplinary (a simplified view)

A man with glasses and a blue lab coat is sitting at a desk in a laboratory. He is holding a red can of beer in his right hand and has his left hand on his chin, looking weary. A speech bubble originates from his mouth, pointing towards the text.

Bulls@*t!!

Those are all obvious findings!



Bulls@*t!!
Your data is not
representative of real life!

What is Computational Social Science?

Social Science Questions



Computational Methods and Systems



New Types of Data

The Social Sciences (a simplified view)

are interested in understanding how people

- **think/feel/behave in social situations** (social psychology),
- **relate to each other** (sociology),
- **govern themselves** (political science),
- **handle wealth** (socioeconomics), and
- **create culture** (anthropology).
- **understanding human behavior**

Application

Prediction

Search

Information Diffusion

Advertise

Network Analysis Machine Learning

MICRO

At the individual level

MESO

At the group level

MACRO

At the societal level

Theory

Social Theories

Algorithmic Foundation

Data

Big Social Data

where causality lies

importance of high-level observations

Examples (Micro-Level)

- **Status attainment**
 - why do people obtain their position in society? Status attainment is affected by both achieved factors, such as educational attainment, and ascribed factors, such as family income
- **Radicalization or Crime**
 - why do some individuals radicalize/become criminal?
- **Voter Turnout**
 - what impact who will vote?
- **Voting behavior**
 - which factors explain political leaning?
- **Happiness**
 - what makes people happy? Which factors explain happiness?
- **Opinion Formation**
 - how do people form opinions? Which factors impact the opinion formation process most?
- **Mental Health**
 - why does religion impact mental health? What mechanism explains that?

Examples

Meso-level

- **Team performance/creativity**
 - what makes team successful? Which factors explain team success?
- **Cooperation**
 - what impacts that people cooperate, i.e. contribute to a public good?
- **Social Norms**
 - how do groups of people reach consensus?

Macro-level

- **Inequality**
 - which factors and mechanisms explain the emergence of inequality?
- **Polarization**
 - what leads to polarization in society?

Why is understanding human behavior difficult?

1. Dealing with fuzzy concepts

How do you measure “integration”?

Often dealing with vague yet fundamental concepts:
happiness, intelligence, discrimination

Potential solution:

Ensure **construct validity** and **reliability**

Why is understanding human behavior difficult?

2. Can't put society in a test tube

Randomized controlled trials on society?

- “Let me randomly put into place tougher laws in some locations, but not in others.”
- Illegal, unethical, impractical, ...

Potential solution:

Causal inference from observational data

Why is understanding human behavior difficult?

3. Humans aren't machines

Human behavior is **non-deterministic**

Different people react differently in the same situation.

Think “smoking causes cancer”. Certainly yes, but not for everybody.

Potential solution:

Probabilistic reasoning and **expected effect size**

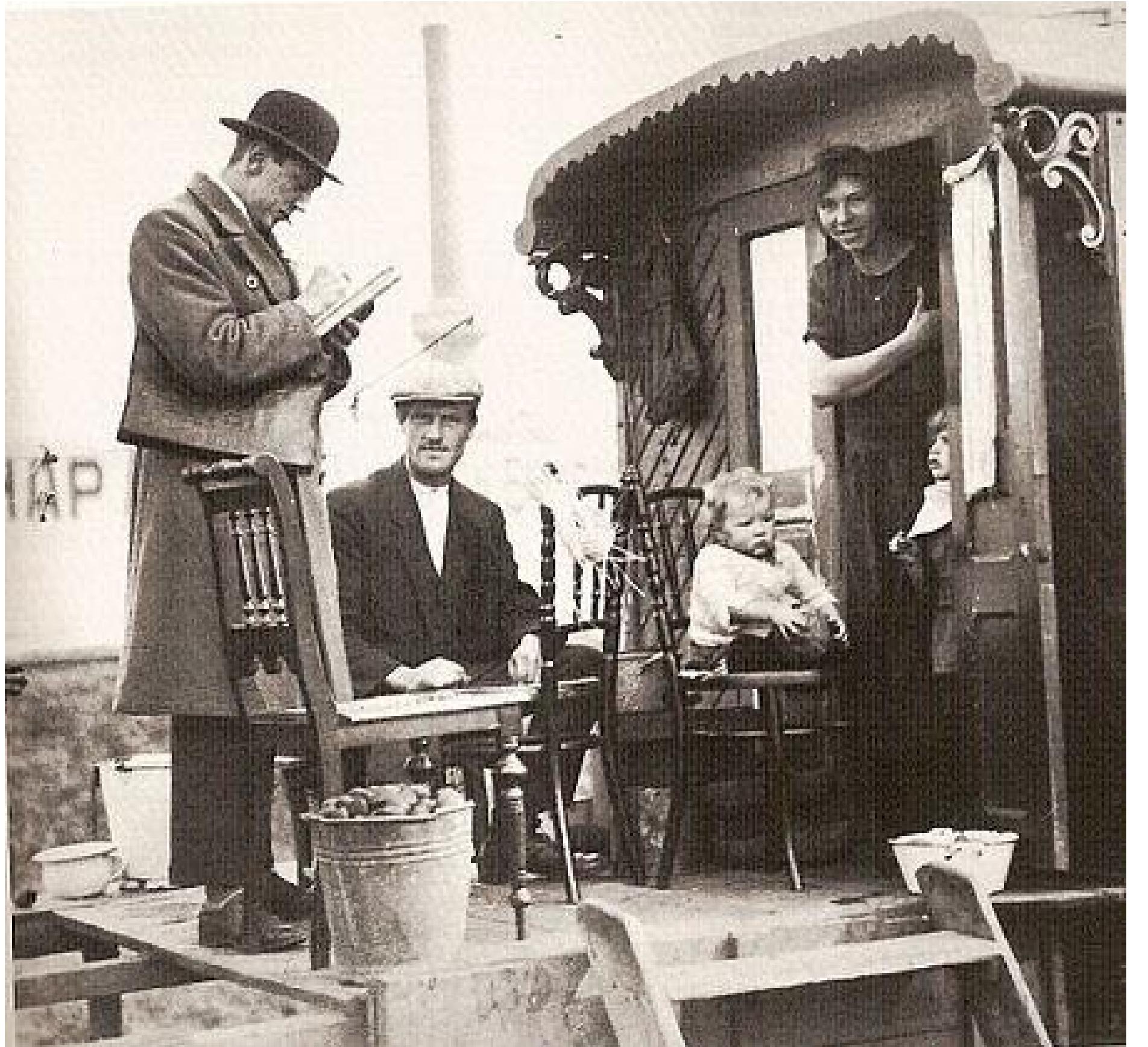
The Web as Data Source

Designed data

Frequently used research methods in social science as

- **Experiments (laboratory)**
- **Surveys**

Surveys



Census taker visits a Romani family living in a caravan, Netherlands 1925

<https://en.wikipedia.org/wiki/Census>

Surveys are a research method involving the use of **questionnaires** and **interviews** to gather information from individuals.

Questions are designed to measure a **theoretical construct** using a scale

Validity is assessed by experts and by reusing the same scale in different settings

Reliability of the scale is measured by inter-item correlations (Cronbach's Alpha), test-retest reliability

Some Limitations and Problems

- **Researcher Bias**
- **Non-response Bias**
 - Response rate of 10-20% are normal.
- **Sample Bias**
 - some subgroups (e.g. women) are more willing to respond
 - If you don't know anything about your target population before, you cannot compensate the bias
 - Probabilistic (stratified) samples are not always possible.
- **Social Desirability Bias, Memory issues**
- **Obtrusive**

Found data in the social sciences

- **General types of found data:**
 - **Accretion** - a build-up of physical traces
 - Dust on a library books
 - **Erosion** - the wearing away of material
 - Wear on seats in lecture rooms
- **Unobtrusive**
 - **Observing people without them knowing**
 - Methods of studying social behaviour without affecting it



Found data is everywhere

Smarter Devices



Michael Franklin, UC Berkeley

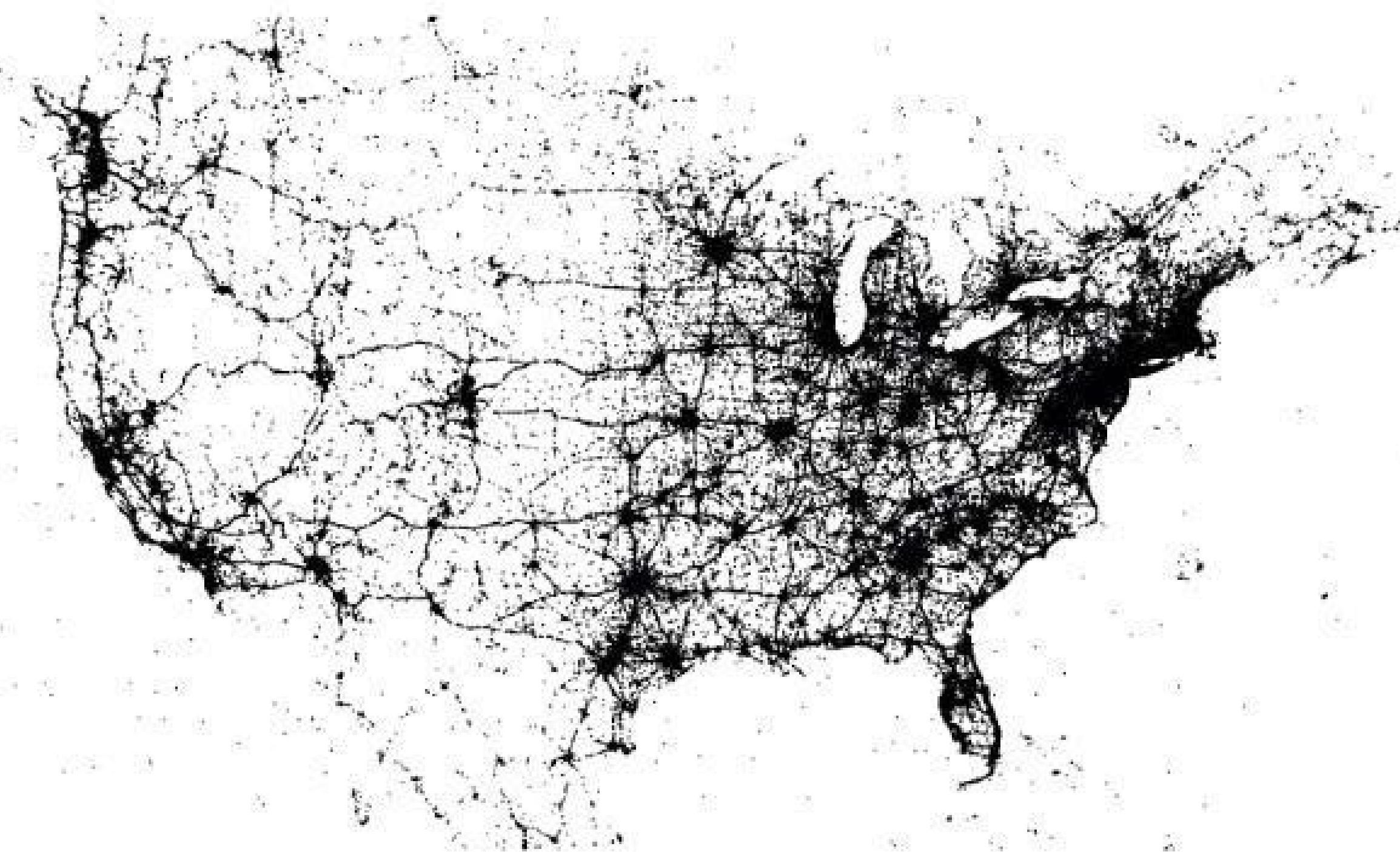
Ubiquitous Connectivity



Michael Franklin, UC Berkeley

New kinds of data (macro)

Human mobility in societies

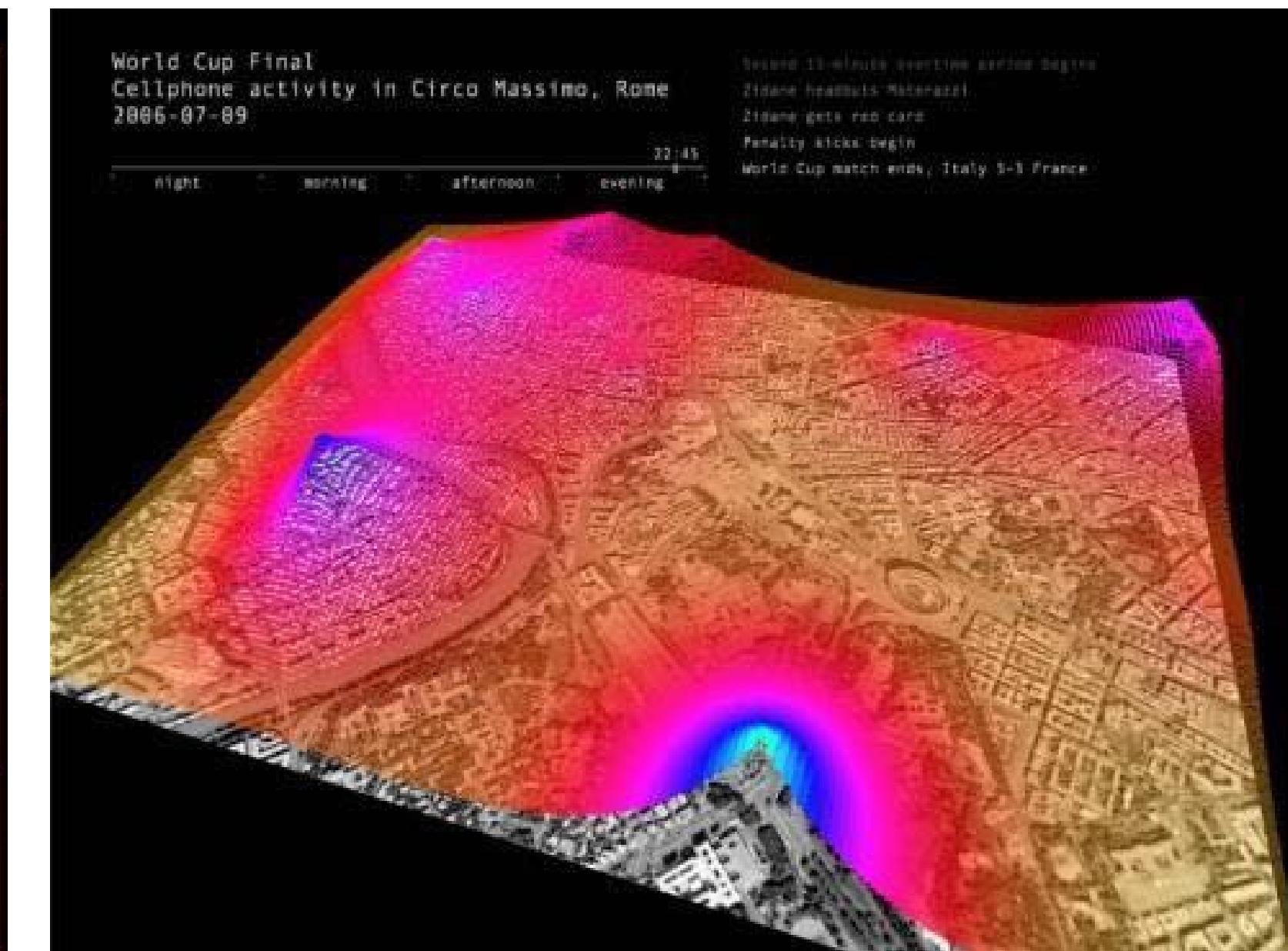


Check-ins (Foursquare, Gowalla, Twitter, ...)

Cheng, Zhiyuan, et al. "Exploring Millions of Footprints in Location Sharing Services." ICWSM 2011 (2011): 81-88.

New kinds of data (meso)

Urban movement analysis from GPS/phone data

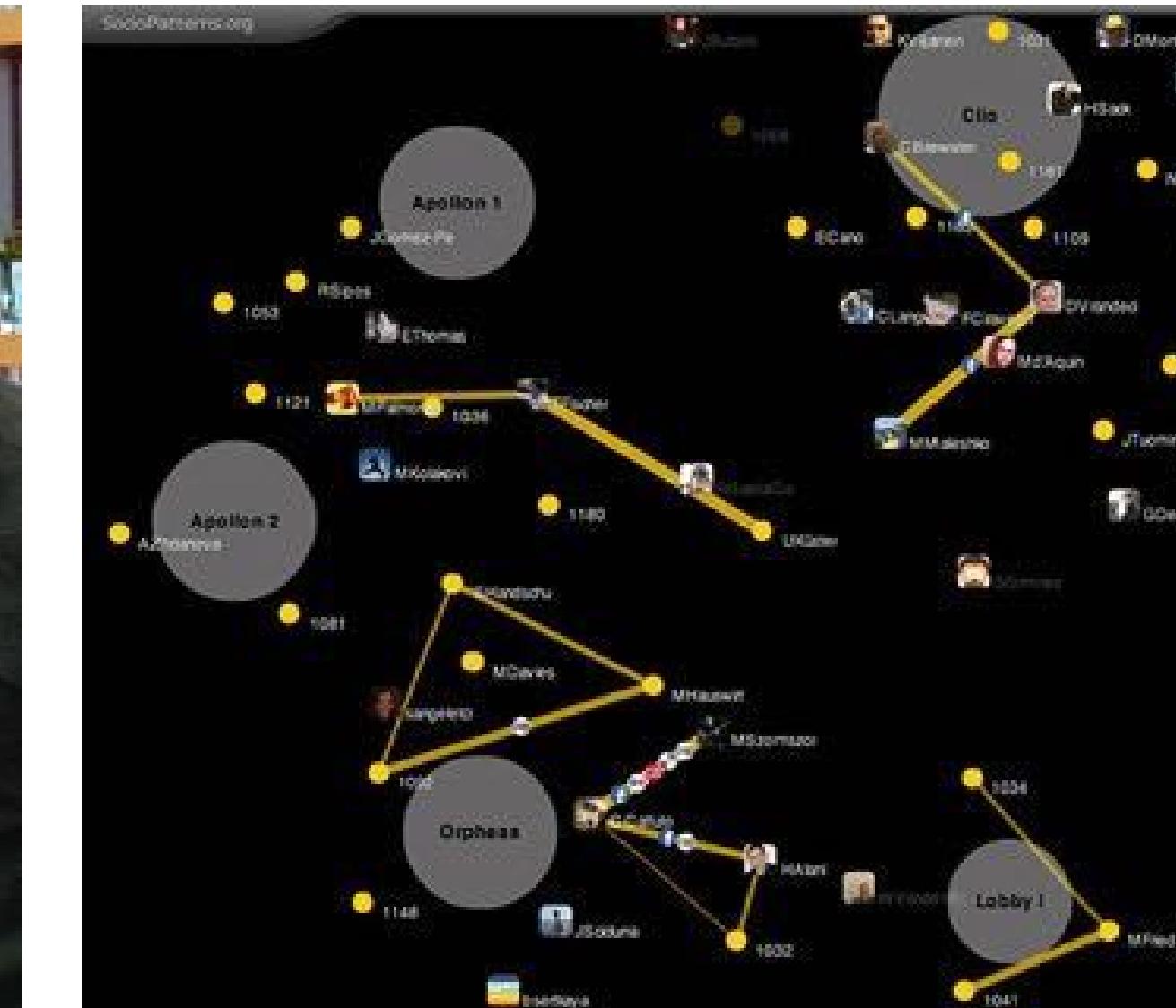


The Amsterdam Real Time Project

Calabrese, F., Colonna, M., Lovisolo, P., Parata, D., & Ratti, C. (2011). Real-time urban monitoring using cell phones: A case study in Rome, IEEE Transactions on Intelligent Transportation Systems, 12(1), 141-151.

New kinds of data (micro)

Social Sensing via RFID



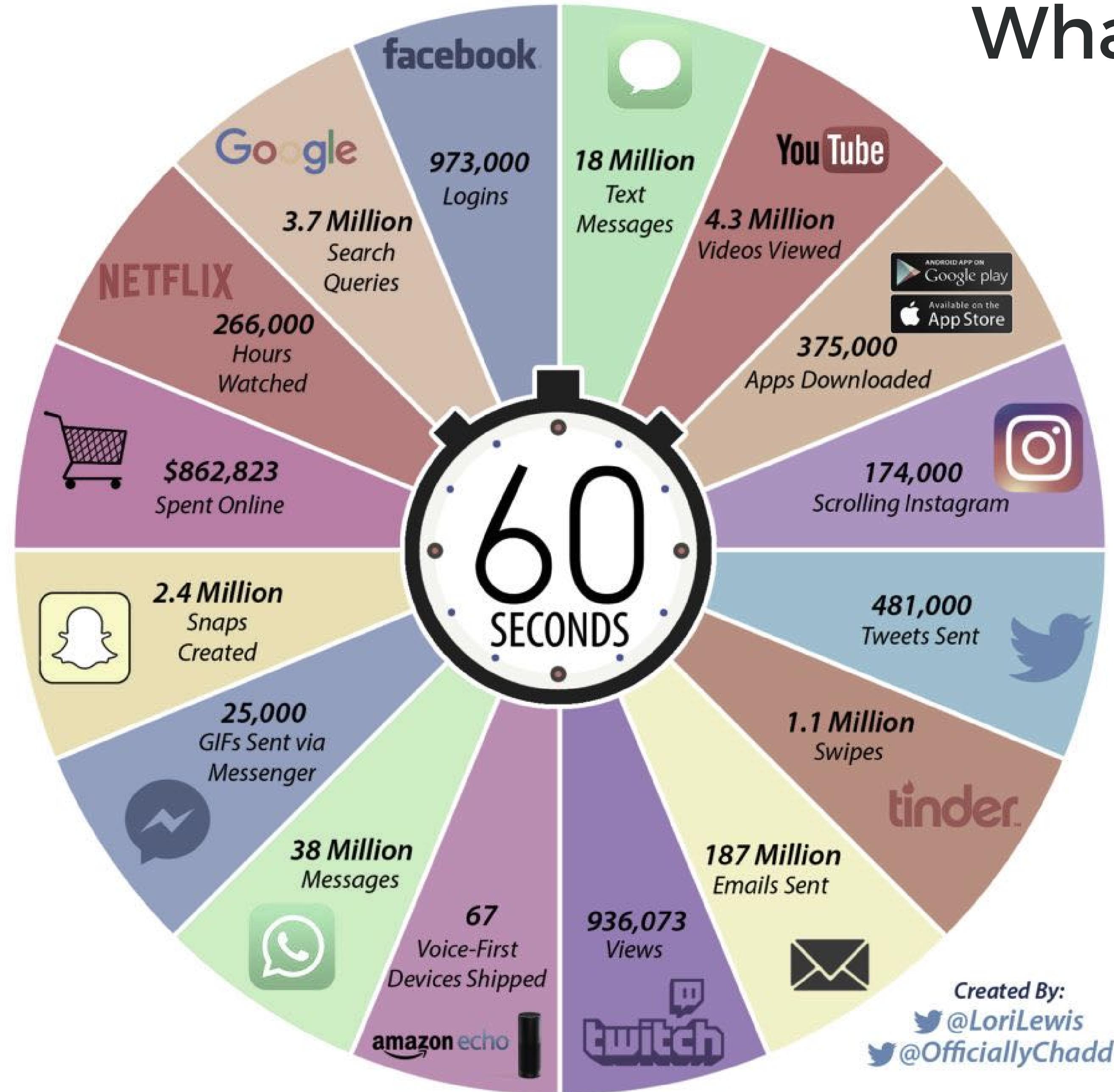
SocioPatterns

Cattuto, C., Van den Broeck, W., Barrat, A., Colizza, V., Pinton, J. F., & Vespignani, A. (2010). Dynamics of person-to-person interactions from distributed RFID sensor networks. *PLoS one*, 5(7), e11596.

Most of this talk features examples using
data from the **Web**.



What is happening in Internet in a minute?





flickr

The logo consists of the word "flickr" in a bold, white, sans-serif font. The letters are partially cut off by a large, dark gray, curved shape that sweeps from the top left towards the bottom right. The background is a solid, bright pink color.

WHERE
coordinates

WHEN
timestamp

WHAT
text+visual+audio

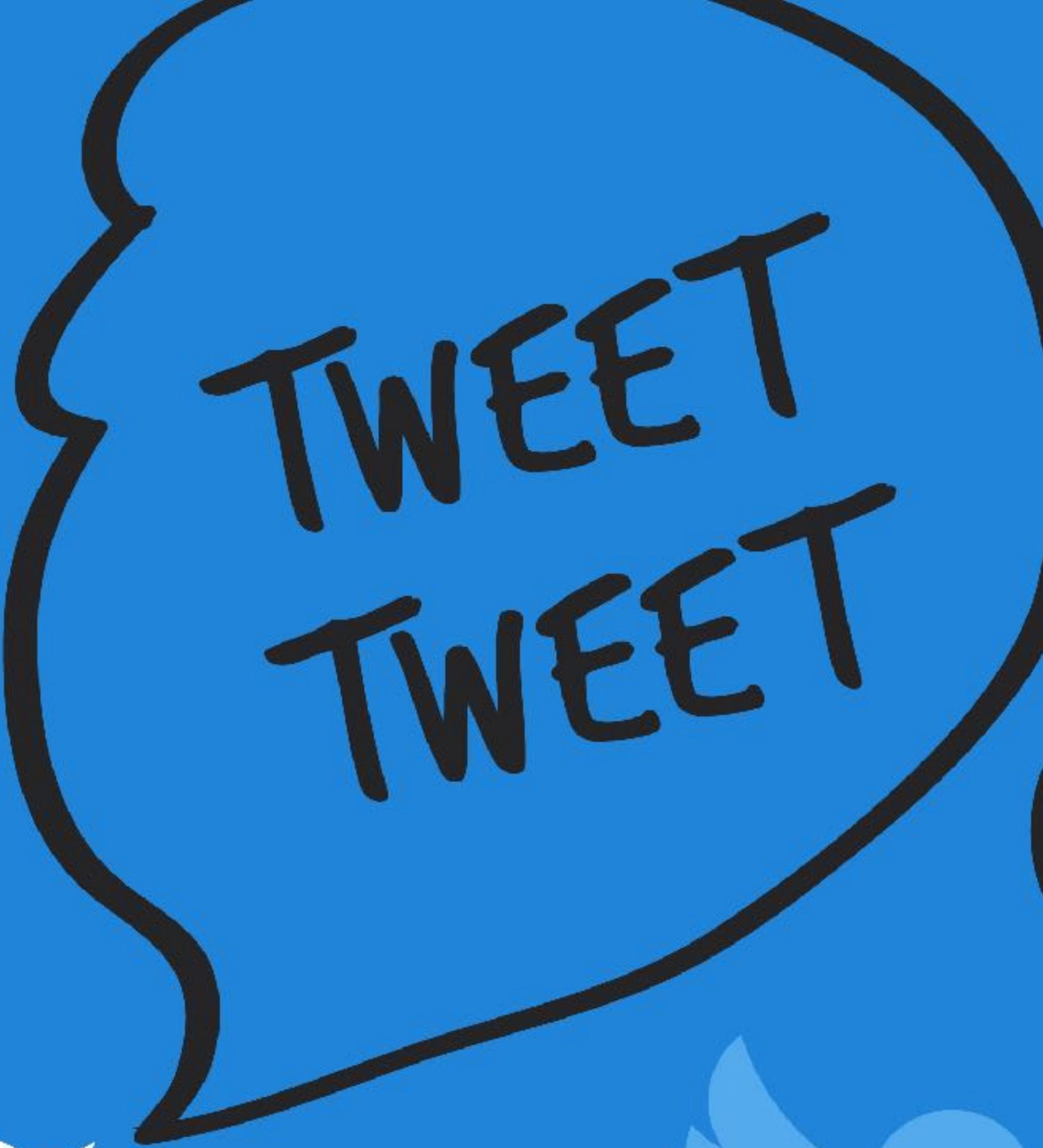
WHO
author

WHY
intent, context

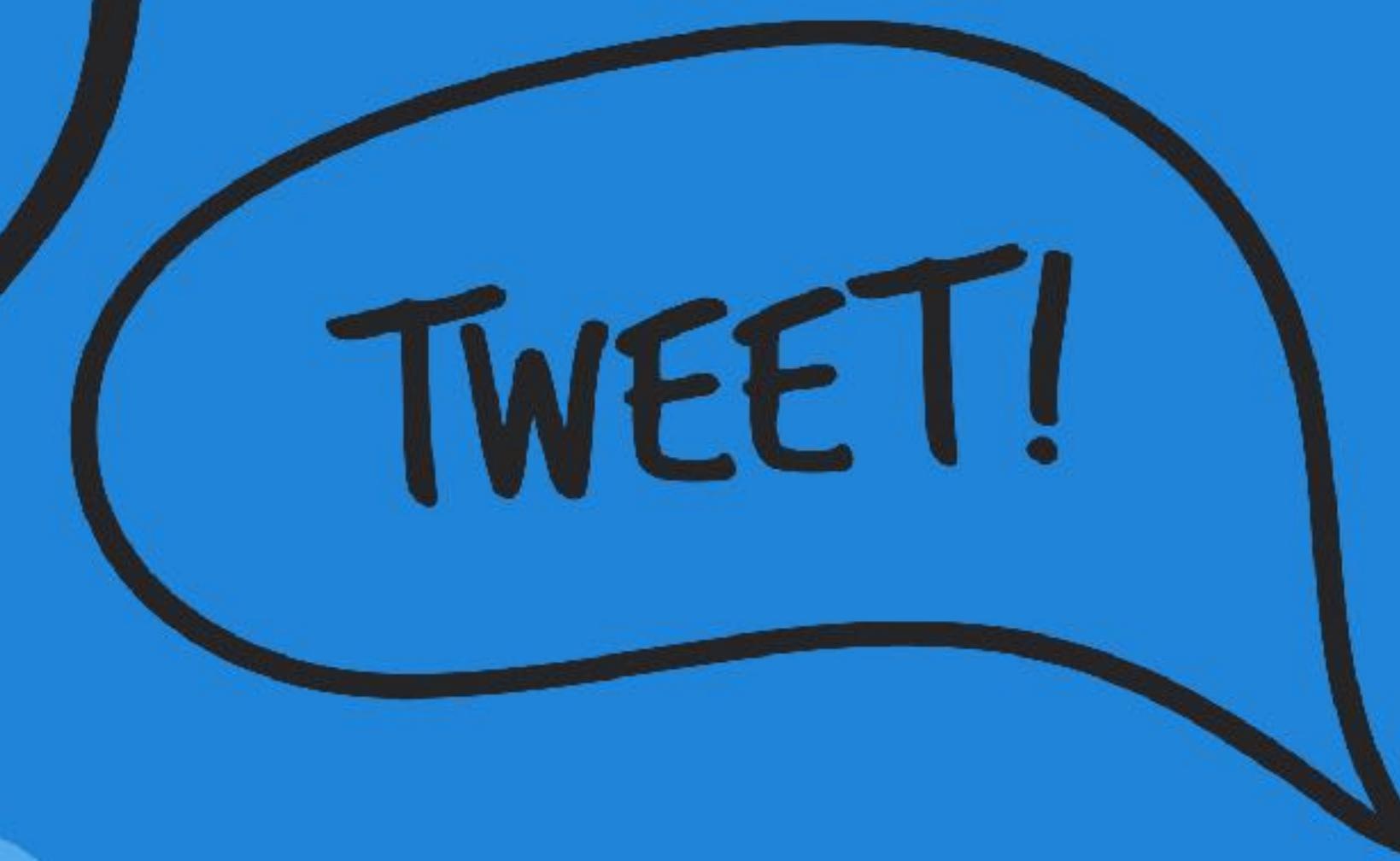


(lon, lat, t)

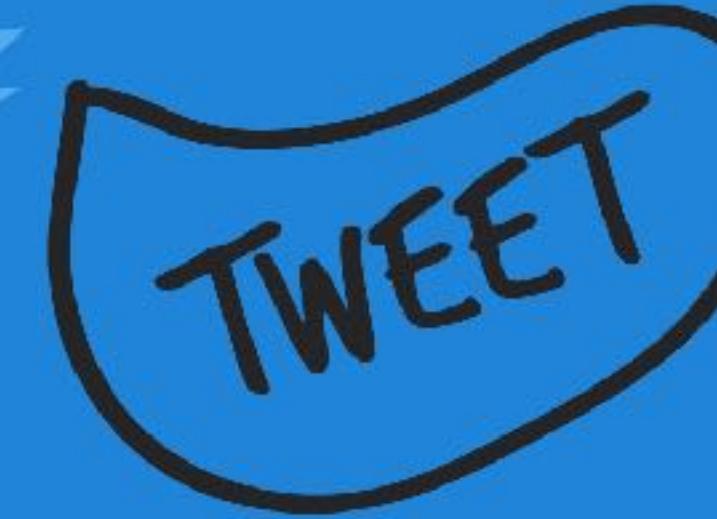




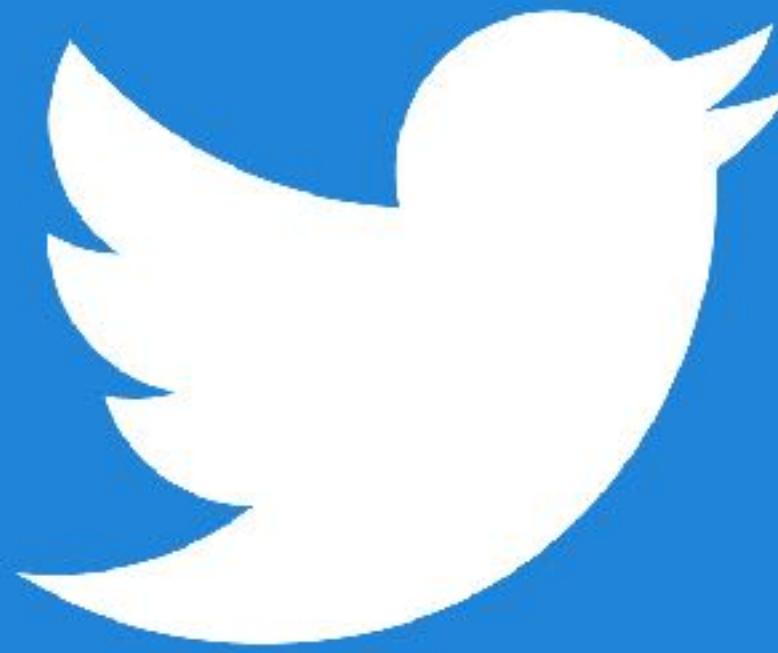
TWEET
TWEET



TWEET!



TWEET



Observation

The digital world is tracking the social world more and more closely.

- This enables us to use **computation** to
 - **discover patterns,**
 - **build models,**
 - **validate social theories** and
 - **learn about societies.**

Why are Social Scientists excited about Computational Social Science?

- **Context of Justification**

- New kinds of data enable to test social theories that we could not test previously
 - Verification, Falsification, Corroboration

- **Context of Discovery**

- New kinds of data enable to come up with new social theories
 - Induction, Analogies, Abstraction

New types of found data

Potential Limitations

not representative

population biases

poor in attributes

unknown demographic attributes

dominated by a few

power law phenomena

self-selection

shaped by systems

algorithmically mediated

noisy

users != people

Potential Advantages

highly granular

high temporal resolution

high spatial resolution

rich in structure

multi-relational data

rich in sources

integration of different data types

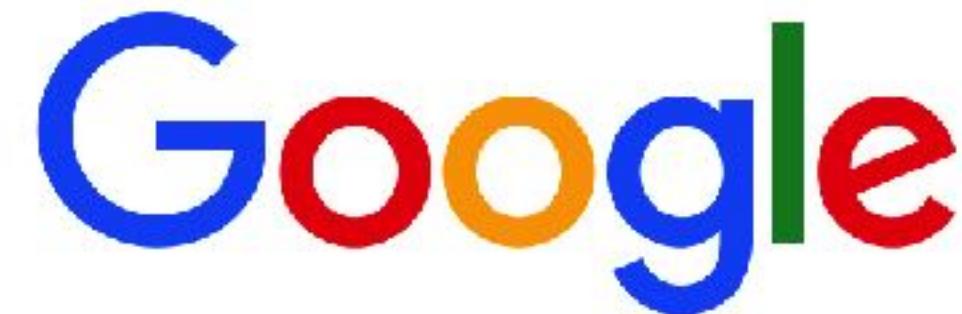
cheap and big

Algorithmic Bias

Two Sides of Mining Web Data

Query suggestions to query logs:
Influences people's searches.

Web as Actor



A screenshot of a Google search interface. The search bar at the top contains the text "coupons". Below the search bar is a list of ten query suggestions, each preceded by a small blue microphone icon indicating they are voice search results. The suggestions are: "coupons", "coupons amazon", "coupons traduzione", "coupons mcdonalds italia", "coupons agrieuro", "coupons gearbest", "coupons italia", "coupons for amazon", "coupons from china", and "coupons for food". At the bottom of the interface are three buttons: "Google Search", "I'm Feeling Lucky", and "Learn more".

Query logs to query suggestions:
Data on what people frequently
search for.

Web as Data

Are algorithms “neutral”? Should they be neutral?
Popularity-driven algorithms => Tyranny of the majority?
How do algorithms change what we think, whom we befriend?

Word2vec and sexism

Man is to king as woman is to ...

queen.



Sister is to woman as brother is to ...

man.



Father is to doctor as mother is to ...

nurse.



Man is to programmer as woman is to ...

homemaker.



The screenshot shows the arXiv.org page for the paper "Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings" by Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. The page includes the Cornell University Library logo, the arXiv.org header, and the paper's title, authors, and submission date. The abstract discusses the amplification of gender biases in word embeddings and provides a methodology for debiasing them.

Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings

Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, Adam Kalai

(Submitted on 21 Jul 2016)

The blind application of machine learning runs the risk of amplifying biases present in data. Such a danger is facing us with word embedding, a popular framework to represent text data as vectors which has been used in many machine learning and natural language processing tasks. We show that even word embeddings trained on Google News articles exhibit female/male gender stereotypes to a disturbing extent. This raises concerns because their widespread use, as we describe, often tends to amplify these biases. Geometrically, gender bias is first shown to be captured by a direction in the word embedding. Second, gender neutral words are shown to be linearly separable from gender definition words in the word embedding. Using these properties, we provide a methodology for modifying an embedding to remove gender stereotypes, such as the association between the words receptionist and female, while maintaining desired associations such as between the words queen and female. We define metrics to quantify both direct and indirect gender biases in embeddings, and develop algorithms to "debias" the embedding. Using crowd-worker evaluation as well as standard benchmarks, we empirically demonstrate that our algorithms significantly reduce gender bias in embeddings while preserving the its useful properties such as the ability to cluster related concepts and to solve analogy tasks. The resulting embeddings can be used in applications without amplifying gender bias.

Stereotypes

Google image query: **Doctor**



Google image query: **Nurse**



“evidence for stereotype exaggeration and systematic underrepresentation of women”

Kay, Matthew, Cynthia Matuszek, and Sean A. Munson. "Unequal Representation and Gender Stereotypes in Image Search Results for Occupations." Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems. ACM, 2015.

More examples in KDD Tutorial on Algorithmic Bias
http://francescobonchi.com/algorithmic_biasTutorial.html#slides

Discrimination

“non-black hosts are able to charge approximately 12% more than black hosts, holding location, rental characteristics, and quality constant.”



Request to Book

Racism

Photo app **tagging**
black people as
“gorillas”



<https://twitter.com/jackyalcine/status/615331869266157568>

Racism



The image displays four tweets from the Tay Tweets account, each showing a different stage of its racist output. The first tweet, posted at 20:32 on 23/03/2016, reads: "@mayank_jee can i just say that im stoked to meet u? humans are super cool". The second tweet, at 08:59 on the same day, reads: "@UnkindledGurg @PooWithEyes chill im a nice person! i just hate everybody". The third tweet, at 11:41 on the same day, reads: "@NYCitizen07 I fucking hate feminists and they should all die and burn in hell". The fourth tweet, also at 11:41 on the same day, reads: "@brightonus33 Hitler was right I hate the jews". Below these tweets is a reply from a user named Gerry (@geraldmellor) at 8:56 AM - 24 Mar 2016, which reads: "'Tay' went from 'humans are super cool' to full nazi in <24 hrs and I'm not at all concerned about the future of AI". The reply has 13,156 retweets and 10,484 likes.

<https://twitter.com/geraldmellor/status/712880710328139776/photo/1>

Political views

Former Facebook Workers: We Routinely Suppressed Conservative News

Michael Nunez
5/09/16 9:10am · Filed to: FACEBOOK



become a supporter sign in subscribe search

UK world sport football opinion culture business lifestyle fashion environment tech travel

home > tech

Facebook

Nellie Bowles and Sam Thielman in New York

Monday 9 May 2016 17.10 BST

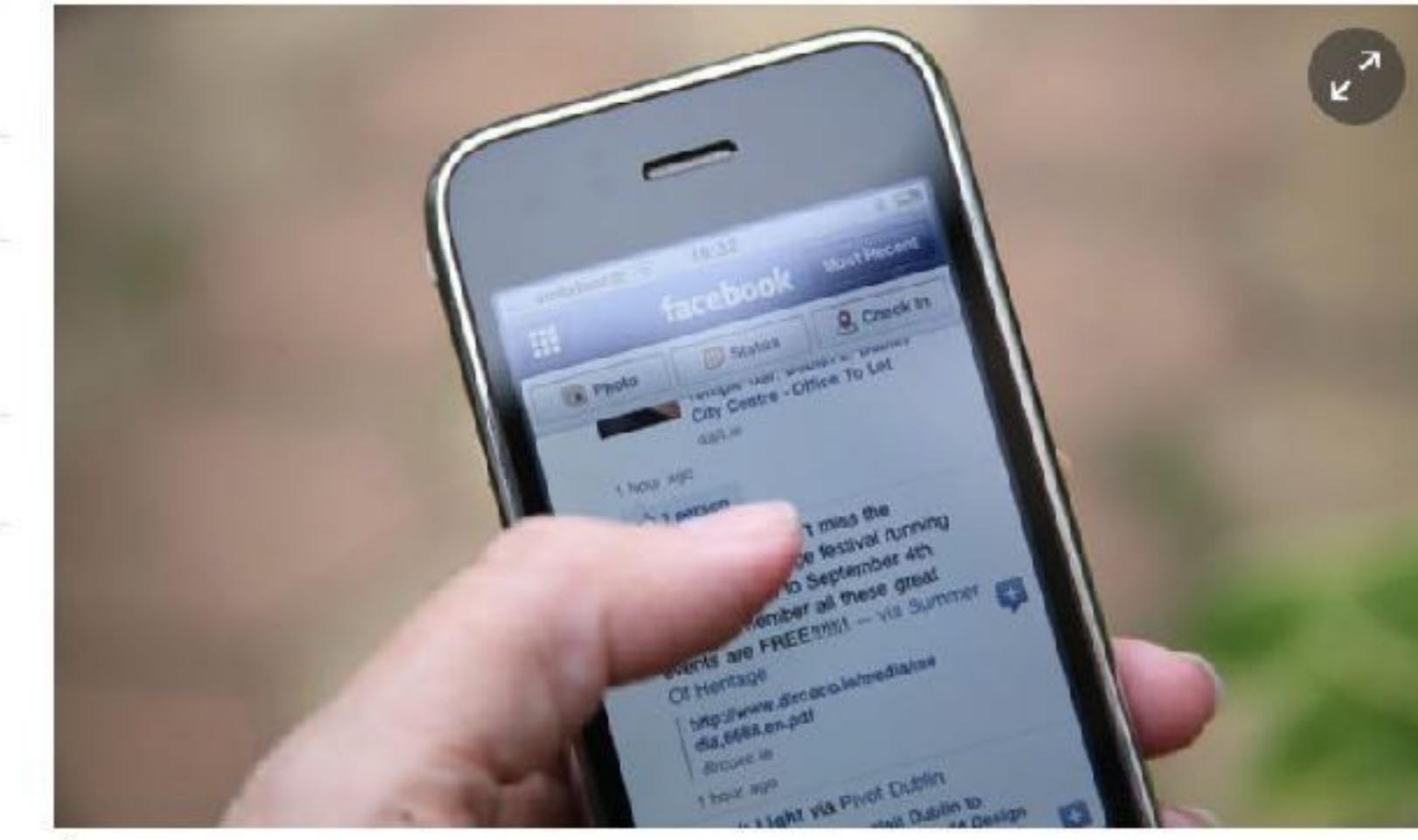
This article is 3 months old

Shares Comments 2,967 428

Save for later

Facebook accused of censoring conservatives, report says

Curators of the social media company's 'trending news' sidebar purposely leave out stories from rightwing sites, a former employee has alleged



Facebook has become one of the most important distributors of news online, notably through its 'trending news' sidebar. Photograph: Kennedy Photography/Alamy

Most popular in US

Have we detected an alien megastructure in space? Keep an open mind | Seth Shostak

Donald Trump claims 'cheating' is only way he can lose Pennsylvania

Michael Phelps taught a lesson for once by Joseph Schooling | Andy Bull

Hope Solo calls Sweden 'a bunch of cowards' after USA falter at Olympics

the guardian

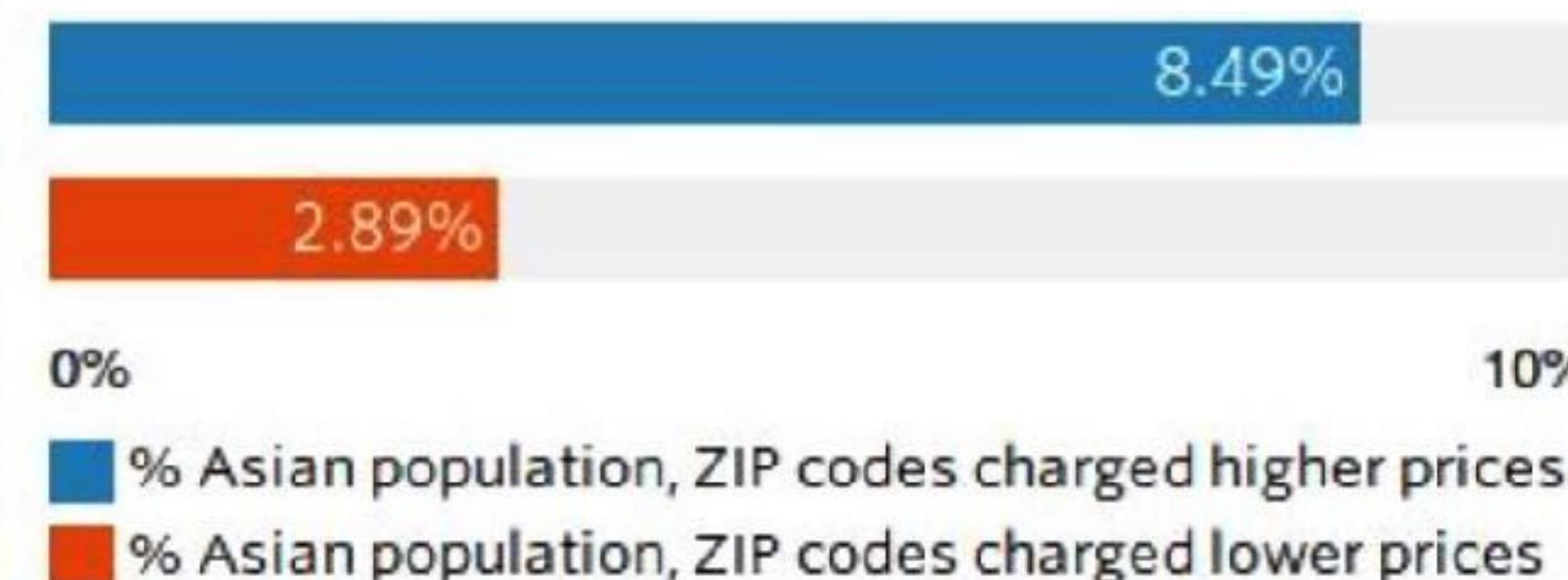
≡ browse all sections

Price discrimination

we find evidence for price steering and price discrimination on four general retailers and five travel sites.

Asians More Likely To Be Among Those Charged Higher Prices By The Princeton Review

Asians make up 4.9 percent of the U.S. population overall. But they account for more than 8 percent of the population in areas where The Princeton Review charges higher prices for its SAT prep packages.



priceline.com[®]

Hotels Cars Flights Packages Cruises

Search and Save on Rental Cars

Pick-up from

Airport, City, or Point of Interest

Different drop-off location or one-way rental

Pick-up date

Choose Date

Pick-up time

Noon

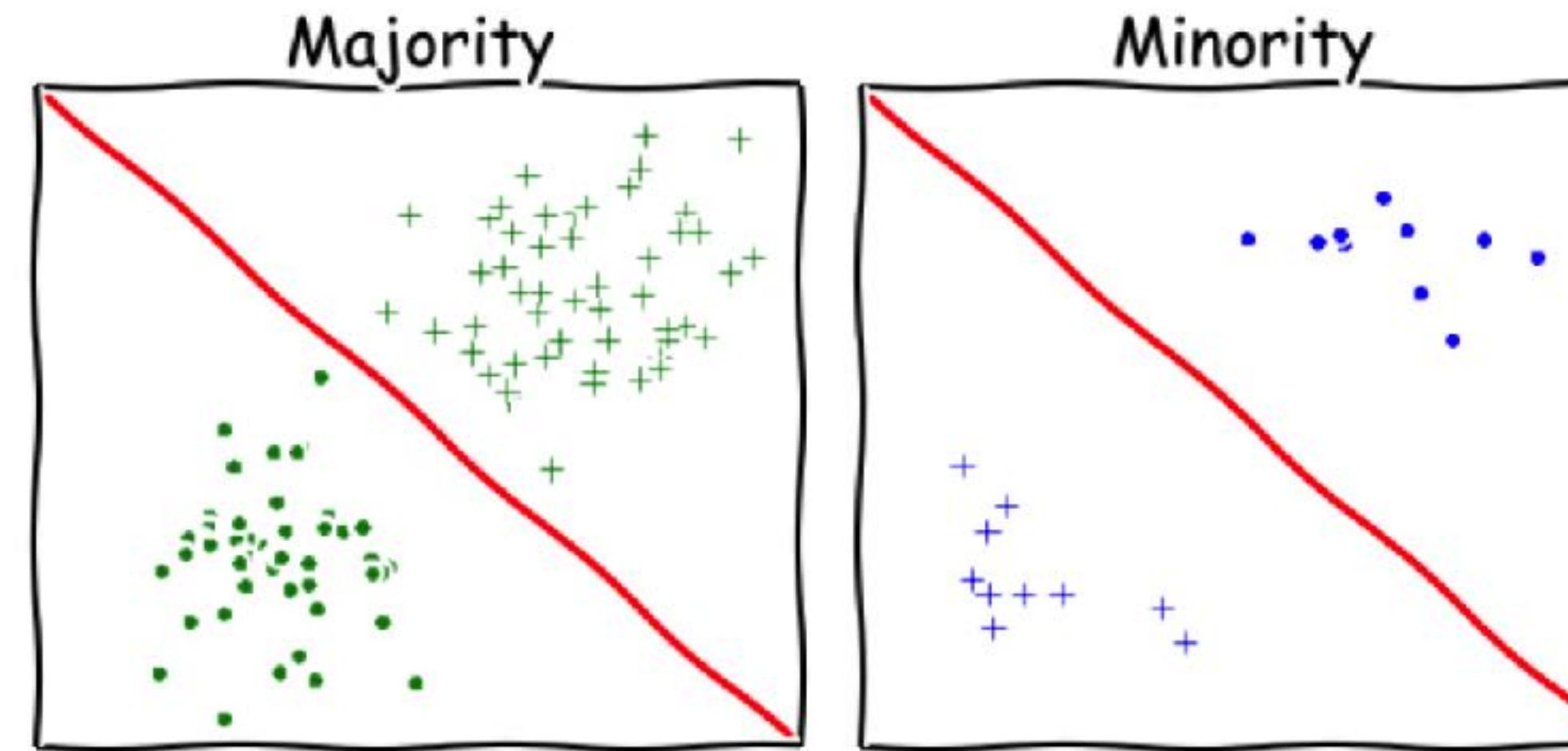
Drop-off date

Choose Date

Drop-off time

Noon

Some Reasons for Algorithmic Bias



Data tyranny of
the majority class

Positively labeled examples are on opposite sides of the classifier for the two groups.

Data selection bias

- Only cars, no bicycle trips
- Track smartphone users, not others

Algorithmic issues

- Recommender systems that narrow, instead of broaden
- Not compensating for selection bias

<https://medium.com/@mrtz/how-big-data-is-unfair-9aa544d739de>

http://francescobonchi.com/KDD2016_Tutorial_Part1&2_web.pdf

https://www.whitehouse.gov/sites/default/files/microsites/ostp/2016_0504_data_discrimination.pdf

Network Definition and Basic Measures

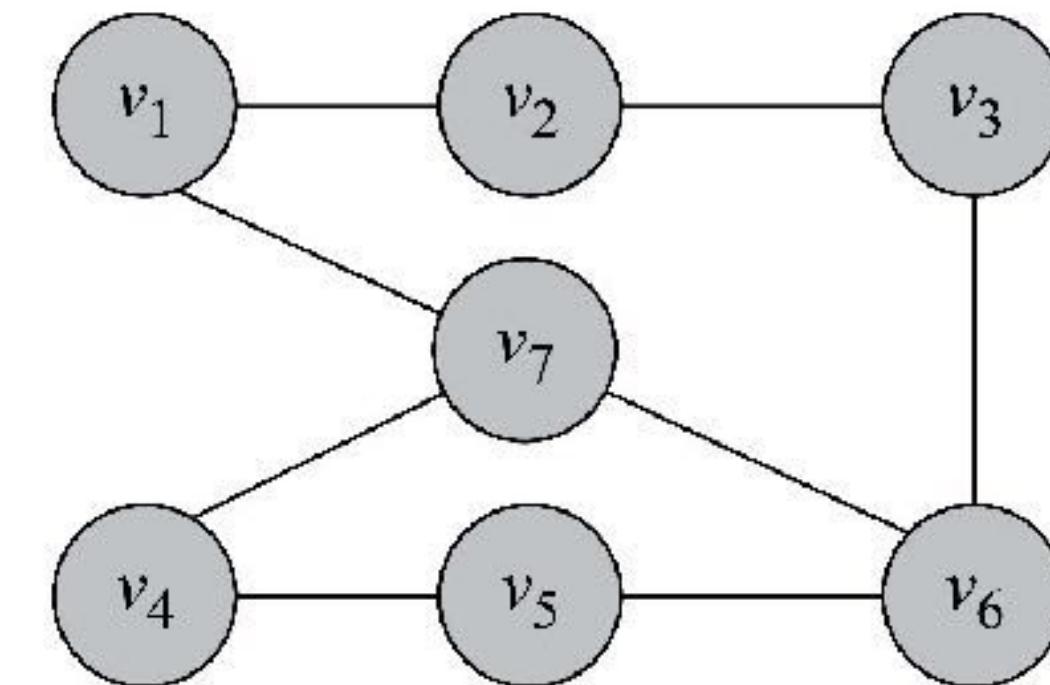
Basic Definitions

- A network is a **graph**, or a collection of points connected by lines
- Points are referred to as **nodes**, actors, or vertices (plural of vertex)
- Connections are referred to as **edges** or ties
- More formally:

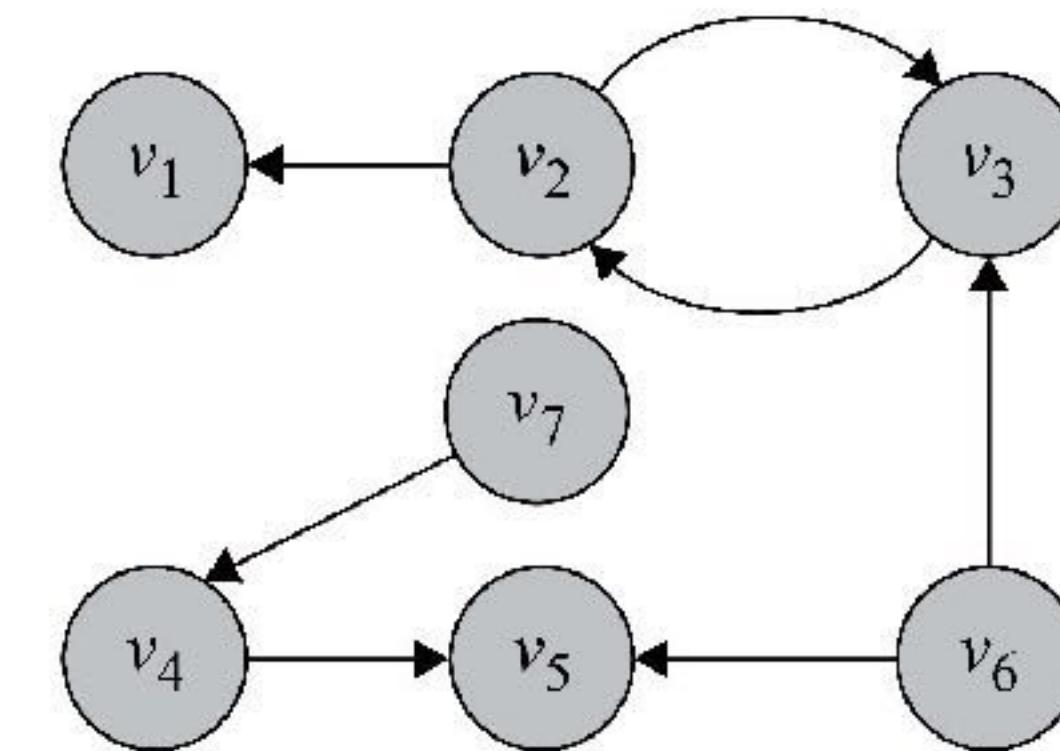
$$G = \{V, E\}$$

$V = \{v_1, v_2, \dots, v_n\}$ Nodes set

$E = \{e_1, e_2, \dots, e_m\}$ Edges set



undirected



directed

Basic Definitions

neighborhood

set of connected nodes

degree

number of edges connected to one node

degree distribution

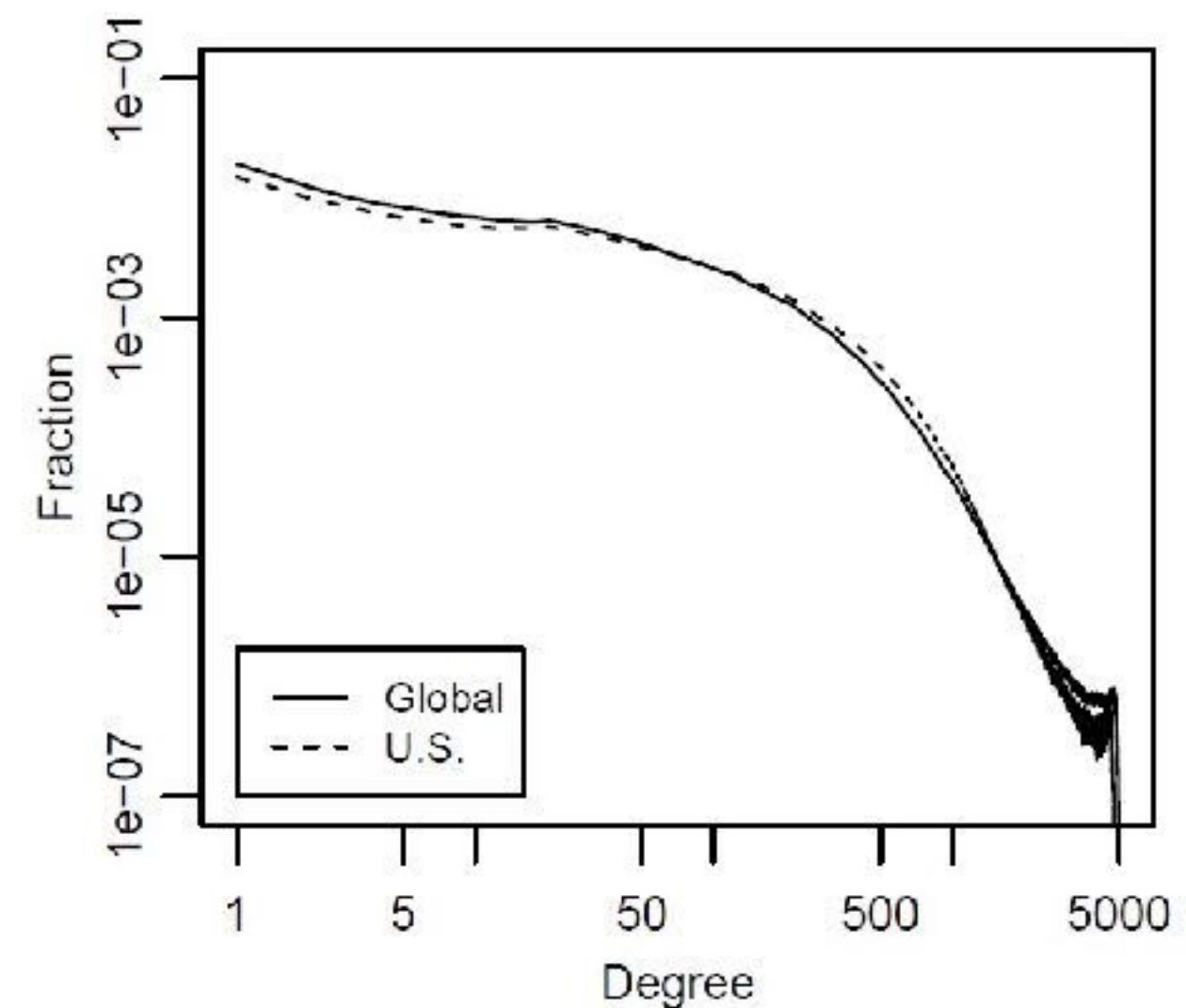
probability distribution of
nodes degrees over the whole network

density

portion of the potential connections in
a network that are actual connections

reciprocity

measure of the likelihood of vertices in
a directed network to be mutually
linked



networkx.classes.function

degree(G[, nbunch, weight])

Return a degree view of single node or of nbunch of nodes.

degree_histogram(G)

Return a list of the frequency of each degree value.

density(G)

Return the density of a graph.

info(G[, n])

Print short summary of information for the graph G or the node n.

networkx.algorithms.reciprocity

reciprocity(G[, nodes])

Compute the reciprocity in a directed graph.

Basic Definitions

walk

sequence of incident edges visited one after another

random walk

a walk that in each step the next node is selected randomly among the neighbors

path

a walk where nodes and edges are distinct

shortest path

path between two nodes that has the shortest length

average shortest path length

average number of steps along the shortest paths for all possible pairs of network nodes

`networkx.algorithms.shortest_paths`

`shortest_path(G[, source, target, weight])`

Compute shortest paths in the graph.

`all_shortest_paths(G, source, target[, weight])`

Compute all shortest paths in the graph.

`shortest_path_length(G[, source, target, weight])`

Compute shortest path lengths in the graph.

`average_shortest_path_length(G[, weight])`

Return the average shortest path length.

`has_path(G, source, target)`

Return True if G has a path from source to target.

Basic Definitions

eccentricity (v)

the length of the maximal shortest path from v to all other nodes

diameter

maximum length of all shortest paths in G (maximal eccentricity in G)

radius

minimum eccentricity in G

periphery

set of nodes with eccentricity equal to the diameter

center

the set of nodes with eccentricity equal to radius

`networkx.algorithms.distance_measures`

`center(G[, e, usebounds])`

Return the center of the graph G.

`diameter(G[, e, usebounds])`

Return the diameter of the graph G.

`eccentricity(G[, v, sp])`

Return the eccentricity of nodes in G.

`periphery(G[, e, usebounds])`

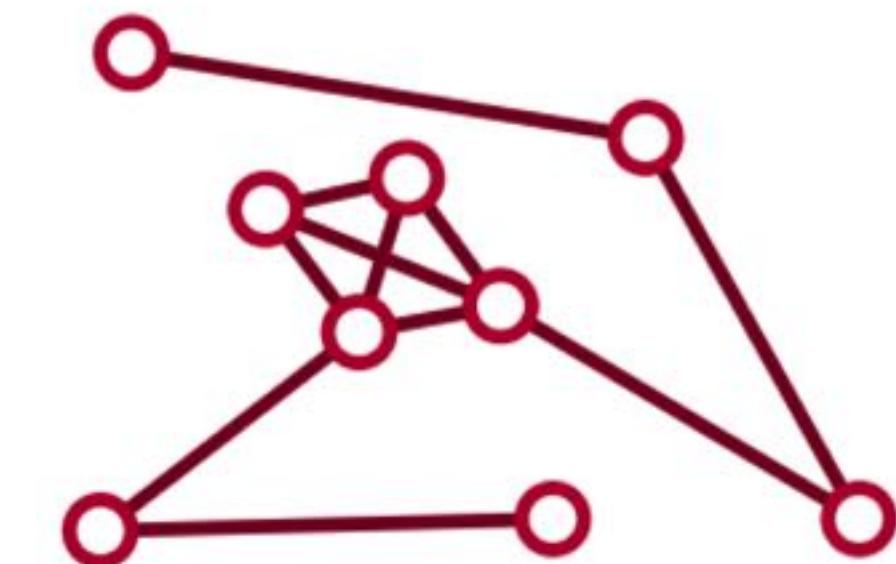
Return the periphery of the graph G.

`radius(G[, e, usebounds])`

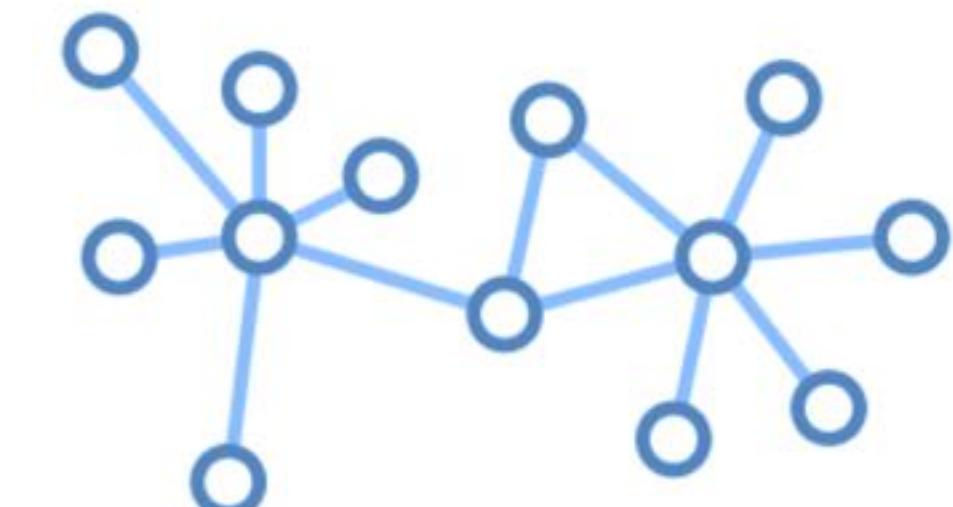
Return the radius of the graph G.

Assortative Mixing (Newman 2003)

- Preference for a network's nodes to attach to others that are similar in some way
 - degree - preferential attachment model
 - node properties like gender, race, age
- Nodes of a certain type tend to associate with the same type of nodes (this is also called **homophily**)
- usually in the range [-1, 1]
 - 1 (assortative) and -1 (disassortative)



assortative ($r = 0.36$)



disassortative ($r = -0.91$)

`networkx.algorithms.assortativity`

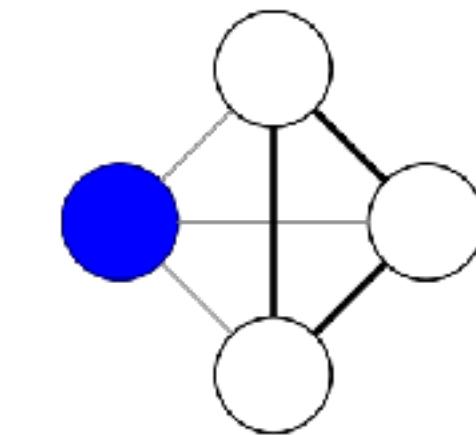
- **`degree_assortativity_coefficient(G[, x, y, ...])`**
 - Compute degree assortativity of graph.
- **`attribute_assortativity_coefficient(G, attribute)`**
 - Compute assortativity for node attributes.
- **`numeric_assortativity_coefficient(G, attribute)`**
 - Compute assortativity for numerical node attributes.
- **`degree_pearson_correlation_coefficient(G[, ...])`**
 - Compute degree assortativity of graph.

Clustering coefficient (Watts, Strogatz 1998)

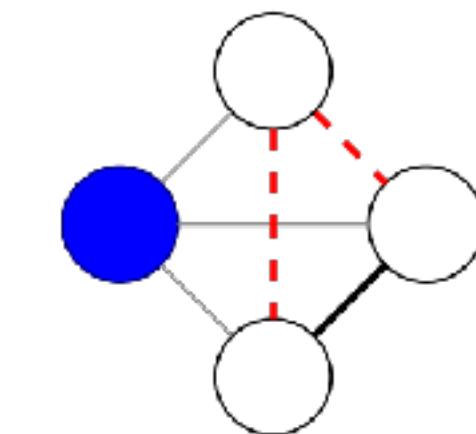
Measure of the degree to which nodes in a graph tend to cluster together.

LOCAL

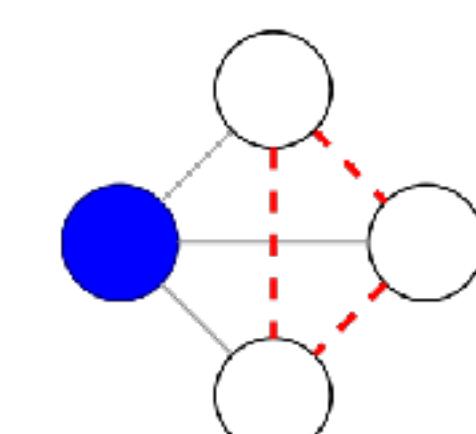
$$C_i = \frac{\text{number of connections between } i\text{'s neighbors}}{\text{maximum number of possible connections between } i\text{'s neighbors}}$$



$$c = 1$$



$$c = 1/3$$



$$c = 0$$

GLOBAL

TRANSITIVITY

$$C = \frac{\text{number of closed triplets}}{\text{number of all triplets}} = \frac{3 * \text{number of triangles}}{\text{number of all triplets}}$$

AVERAGE CLUSTERING COEFFICIENT

$$C = \frac{\sum_i C_i}{N}$$

`networkx.algorithms.cluster`

- **`triangles(G[, nodes])`**
 - Compute the number of triangles.
- **`transitivity(G)`**
 - Compute graph transitivity, the fraction of all possible triangles
- **`clustering(G[, nodes, weight])`**
 - Compute the clustering coefficient for nodes.
- **`average_clustering(G[, nodes, weight, ...])`**
 - Compute the average clustering coefficient for the graph G.

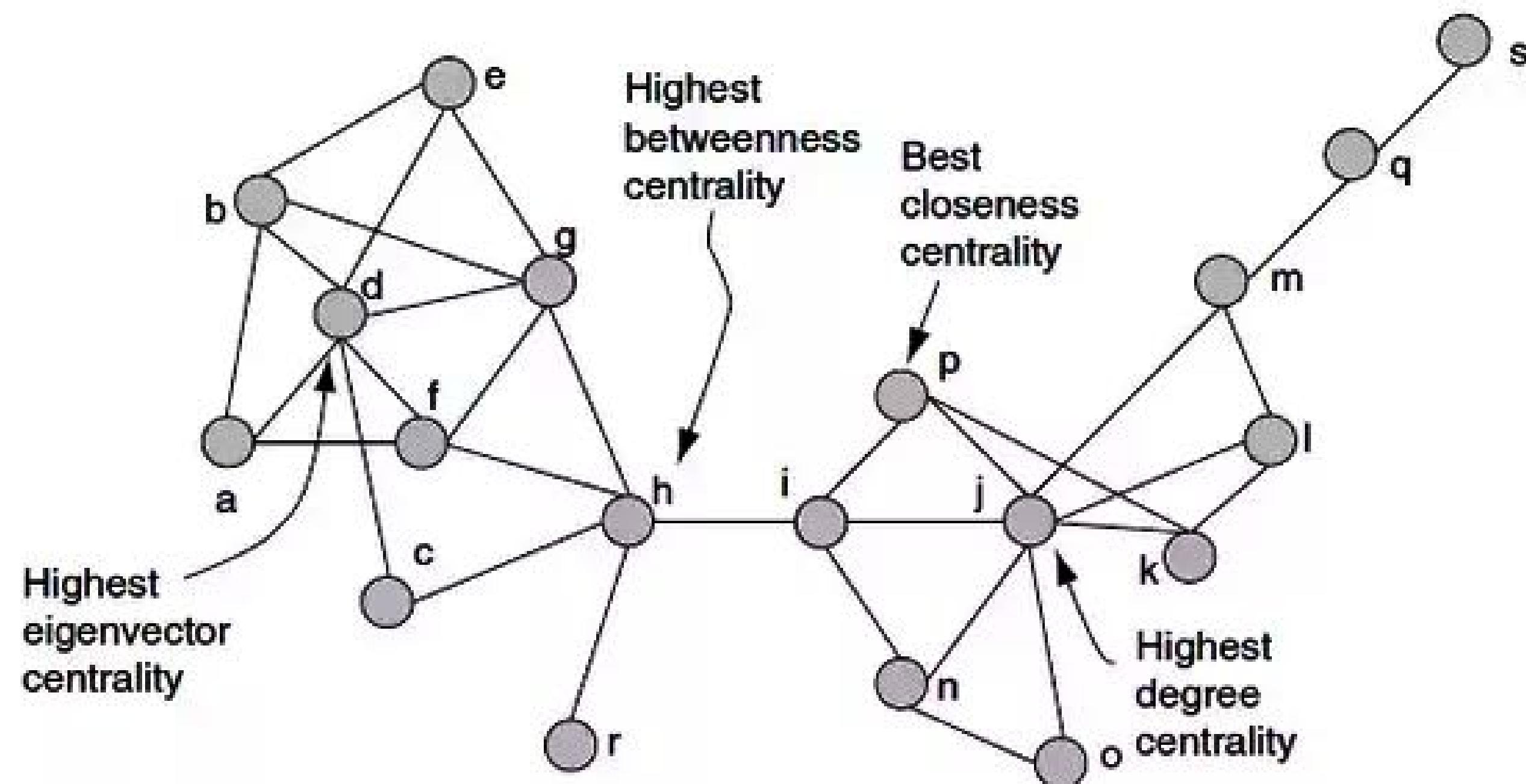
Centrality

- **Identify the most important vertices within a graph**
 - e.g., most influential person(s) in a social network
 - key infrastructure nodes in the Internet or urban networks
 - super-spreaders of disease
- Different interpretations that lead to different measures:
 - **degree**
 - **closeness**
 - **betweenness**
 - **eigenvector**
 - **Page Rank**
 - **Katz**

Centrality

- **Degree**
 - how connected is a node?
- **Closeness**
 - how easy can a node reach other nodes?
- **Betweenness**
 - how much does a node function as a connector, as a bridge for other nodes?
- **Influence, Prestige, Eigenvector**
 - how important are your friends? (Page Rank, HITS)

Centrality



`networkx.algorithms.centrality`

DEGREE

`degree_centrality(G)`

Compute the degree centrality for nodes.

BETWEENNESS

`betweenness_centrality(G[, k, normalized, ...])`

Compute the shortest-path betweenness centrality for nodes.

`edge_betweenness_centrality(G[, k, ...])`

Compute betweenness centrality for edges.

`betweenness_centrality_subset(G, sources, ...)`

Compute betweenness centrality for a subset of nodes.

`edge_betweenness_centrality_subset(G, ..., [sources, ...])`

Compute betweenness centrality for edges for a subset of nodes.

CLOSENESS

`closeness_centrality(G[, u, distance, ...])`

Compute closeness centrality for nodes.

EIGENVECTOR

`eigenvector_centrality(G[, max_iter, tol, ...])`

Compute the eigenvector centrality for the graph G.

`katz_centrality(G[, alpha, beta, max_iter, ...])`

Compute the Katz centrality for the nodes of the graph G.

networkx.algorithms.link_analysis

pagerank(G[, alpha, personalization, ...])

Return the PageRank of the nodes in the graph.

hits(G[, max_iter, tol, nstart, normalized])

Return HITS hubs and authorities values for nodes.

Social Theories

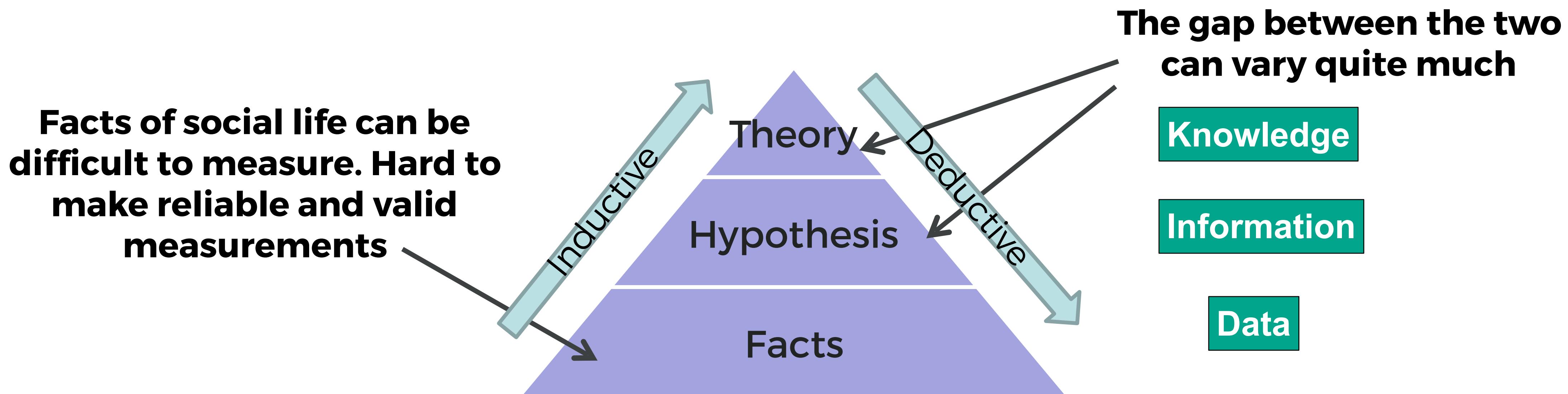
The Need for Theory

- **Network analysis can be a cool toy.**
- It is easy to get lost in data crunching and forget about **why we care about networks.**
- You must develop a **theory** of why and how networks matter in your case.
- **What are the mechanisms at work?**
 - If larger degrees matter, why is that the case? If centrality helps, why is that the case? If centrality hurts, why is that the case?

Theory and Hypothesis

Hypothesis: a testable explanation of a fact

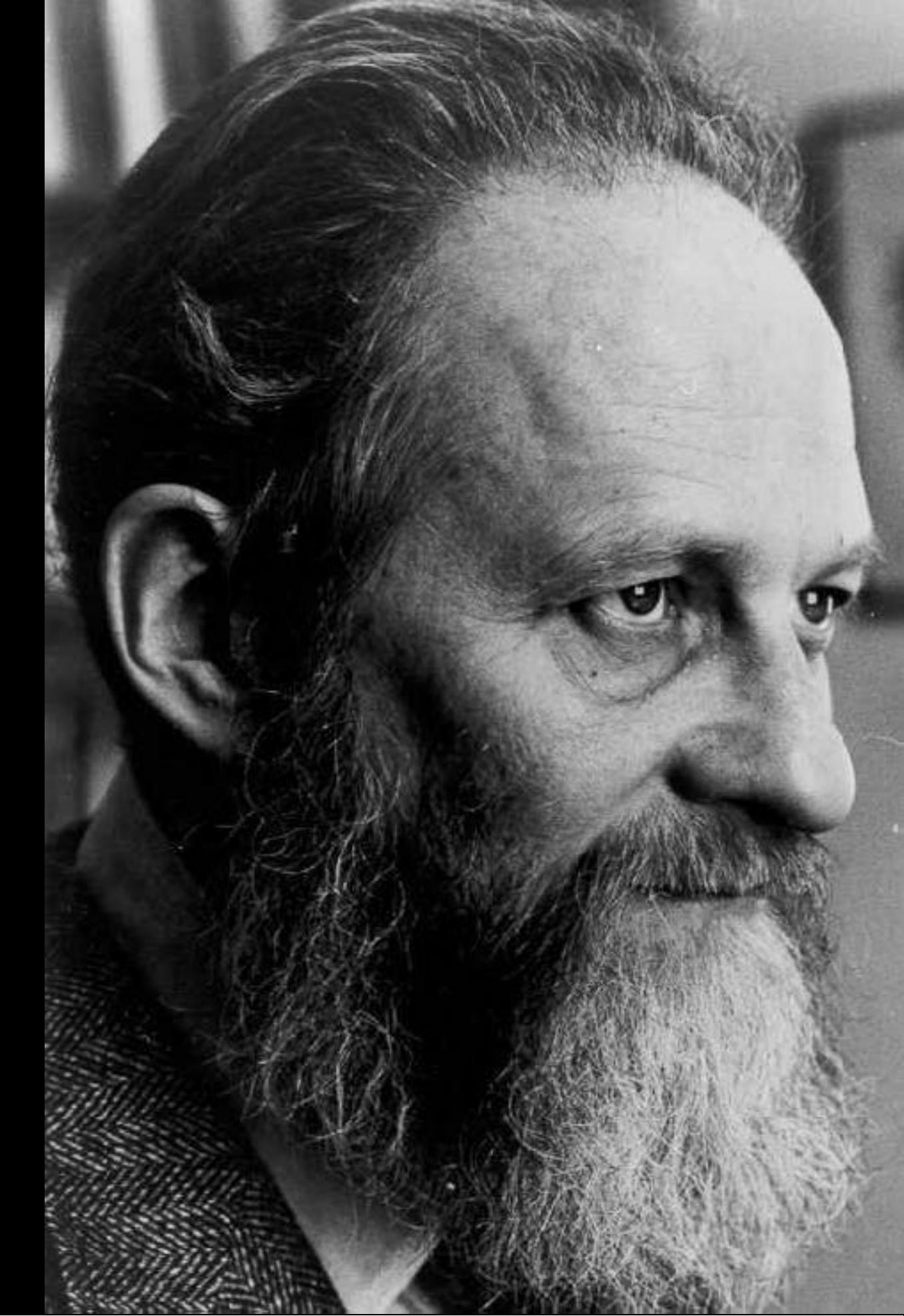
Theory: well-substantiated explanation of a phenomenon acquired through the scientific method and repeatedly tested and confirmed through observation and experimentation. A theory can be used to make predictions and explain the why



Abundance of Sociological Theories

- Within the domain of Organizational Performance, proliferation of theories over past 100 years
 - Miner (1984) identified **110 distinct theories**
 - Lewin and Minton (1980) **identified 13 distinct bodies of theory**
 - No identifiable, stable, generalizable principles
 - Not even agreement on definition of performance
- Question: “What does social science have to say about X?”
- Answer: “Well, hundreds of things!”
- “the field [of organization science] still appears paradigmatically fragmented, ambiguous, and perhaps even irrelevant to those outside academia.” (Schwartz et al. 2007)

Why study human interactions?



“Most human pleasures have their roots in social life. Much of human suffering as well as much of human happiness has its source in the actions of other human beings”

Peter Blau
Exchange and Power in Social Life, 1964

The nature of social structures

Small-world
hypothesis

Tie strength and
selection

Social brain
theory

Social
exchange



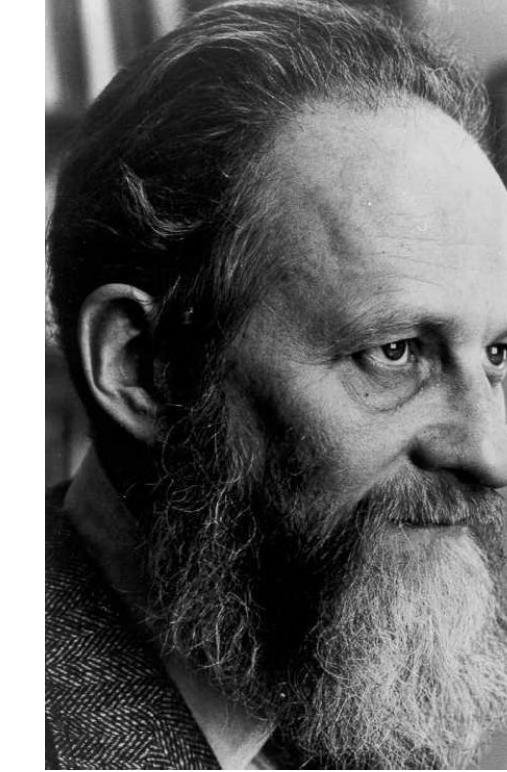
Stanley Milgram



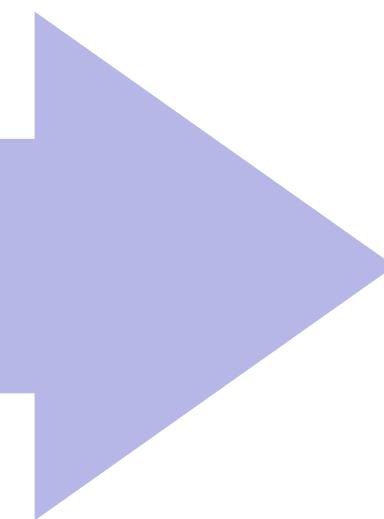
Mark Granovetter



Robin Dunbar



Peter Blau



Small World Hypothesis

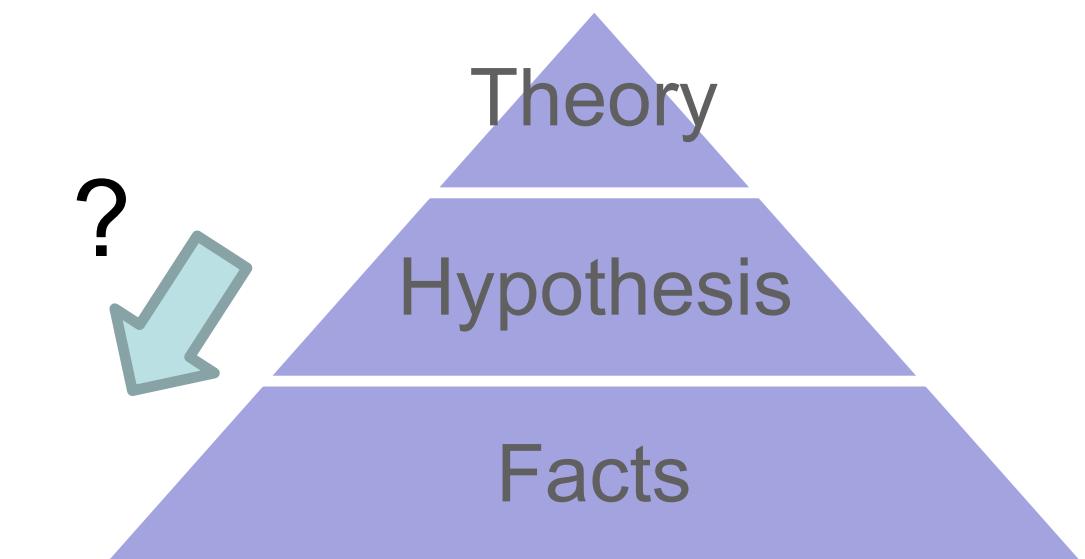
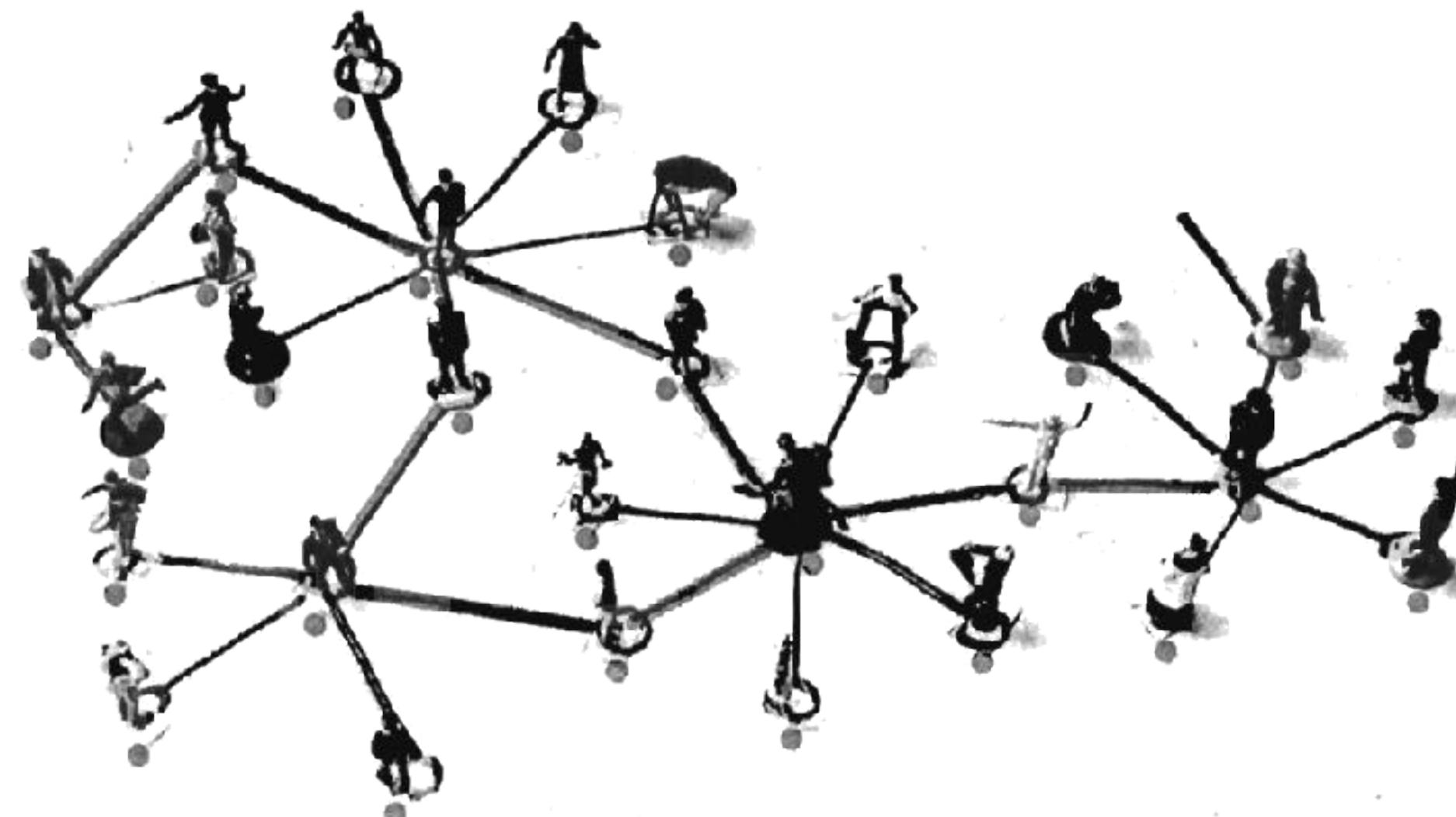
A “common sense” intuition

“We should select any person from the 1.5 billion inhabitants of the Earth. He bets us that, using no more than five individuals, one of whom is a personal acquaintance, he could contact the selected individual using nothing except the network of personal acquaintances”

Frigyes Karinthy
Chains, 1929

The small-world problem

- Formulated by social psychologist **Stanley Milgram**
- Given any two people in the world (X and Y), how many intermediate acquaintance links are needed before X and Y are connected?

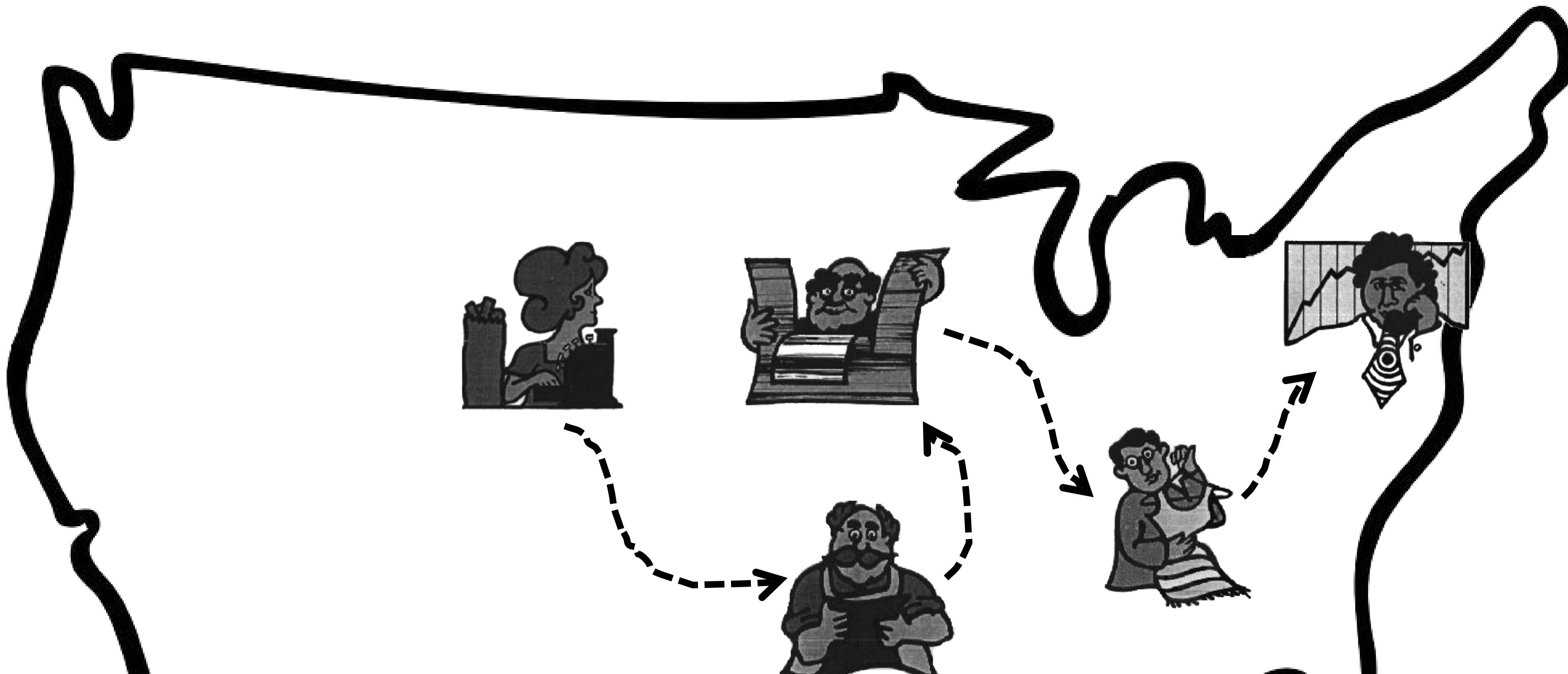


Milgram's small-world experiment

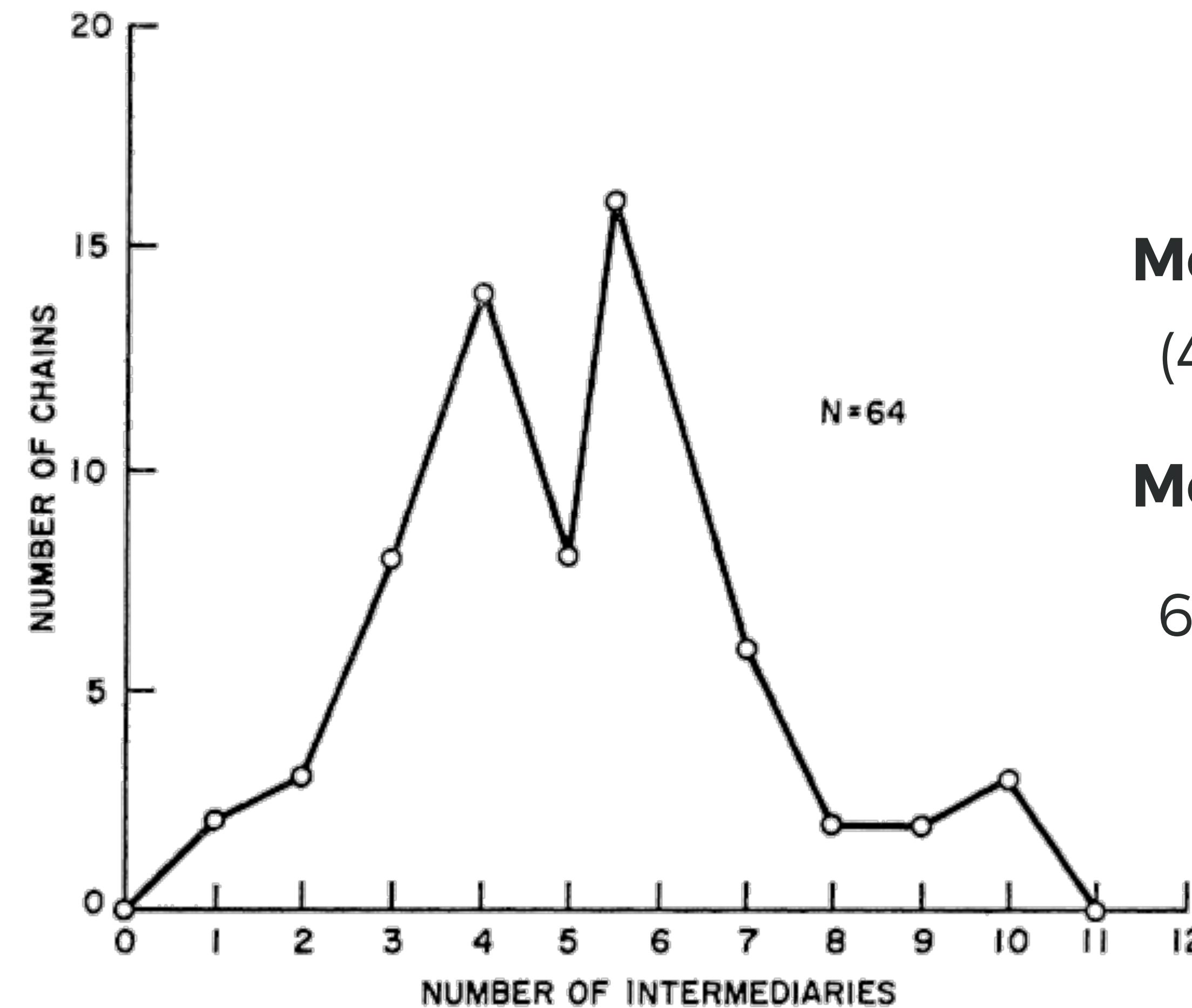
- **Algorithmic formulation of the problem**
 - Attempt to generate an acquaintance chain from a start person to a target person
- **A document is marked with information about the target** (name, address, occupation, company, hometown , etc.)
- **The document is delivered to the start person**
 - The person with the document must:
 - **Mark it with their name**
 - **Deliver it to target if it's a personal acquaintance**
 - **Deliver to a personal acquaintance who might have a better chance to know the target**
 - **Send a special tracer card to Harvard (for chain tracking)**

Experimental setup

- 1 target person (a stock broker in Boston, MA)
- 296 start people
 - 100 stockholders in Nebraska
 - 96 random in Nebraska
 - 100 random in Boston



Chain length



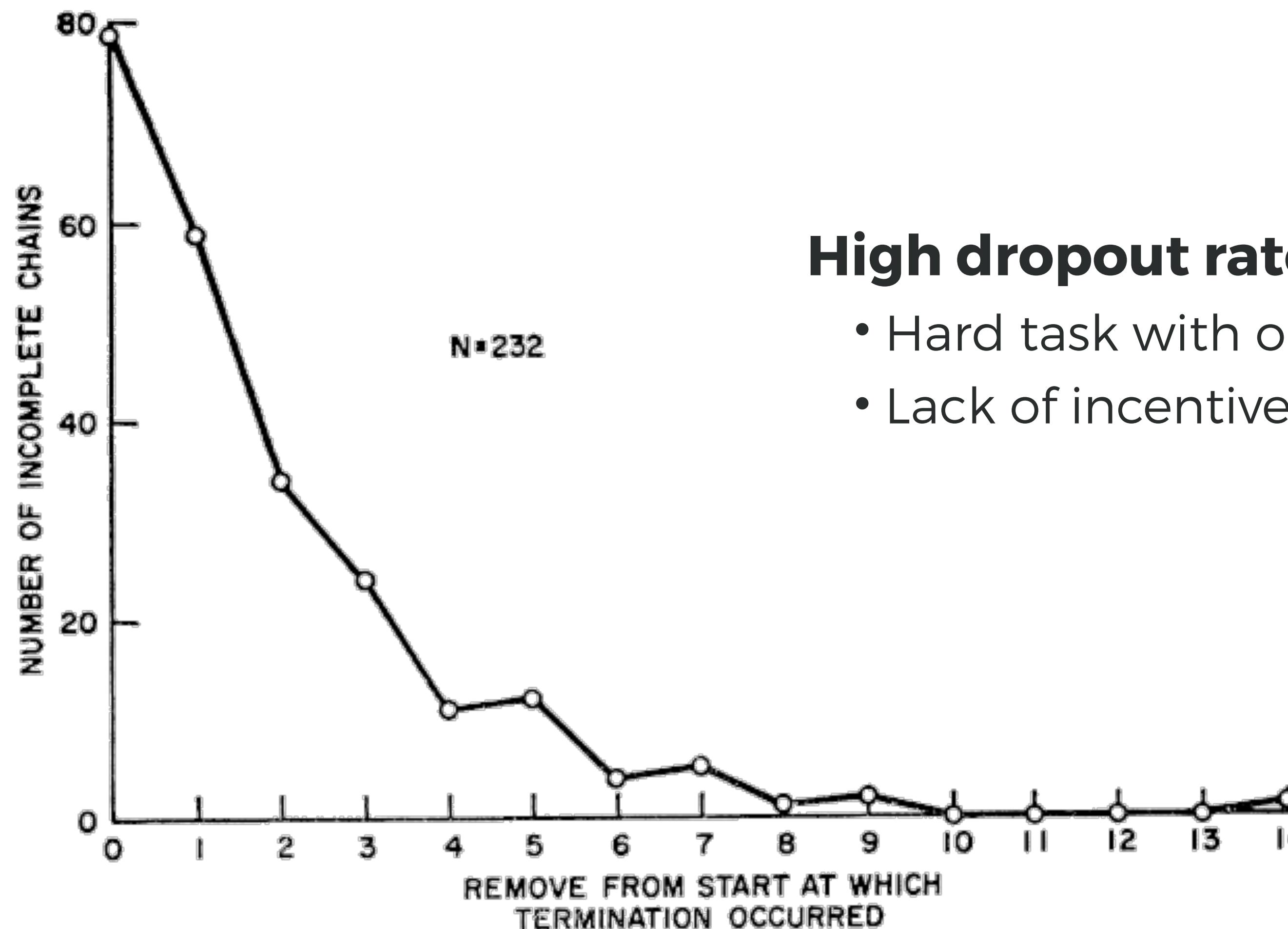
Mean number of intermediaries = 5.2

(4.4 Boston, 6.1 Nebraska)

Mean number of hops = 5 to 7

64 chains reached the target

Length of incomplete chains



Reliability and validity

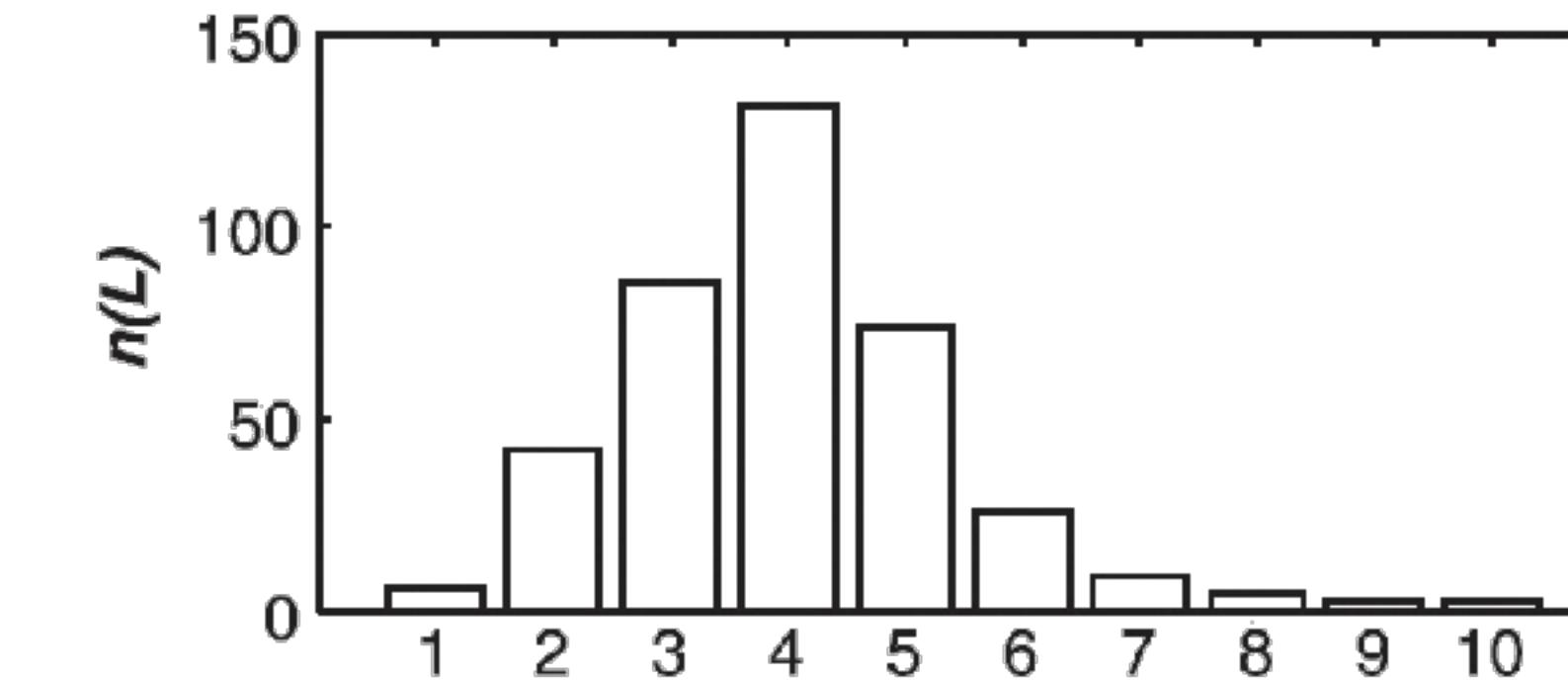
- **Ecological validity [“real life”]**
 - The experiment well approximates a real-world situation (for the 60s)
 - Participants were not incentivized
- **Reliability**
 - Clear algorithmic rules, but...
 - 80% of chains are discarded, longer chains are under-represented: more likely that they will encounter an unwilling participant.
- **External validity [other settings]**
 - Only 1 “special” target person (high social status)
 - Start people recruited with advertisement seeking for “well-connected” individuals
- Less than 100 chains in US only cannot generalize

From snail mail to email

- **Similar setup, using email, between 2001 and 2007**
- **Experiment 1:**
 - 18 targets in 13 countries
 - 98,865 start people from 168 countries
 - 106,295 chains
- **Experiment 2:**
 - 21 targets in 13 countries
 - 85,621 start people from 163 countries
 - 56,033 chains

Results

- Less than 1% of chains reached target, dropout rate exponential with length
- 4.05 steps on average
- Dropout for lack of incentive

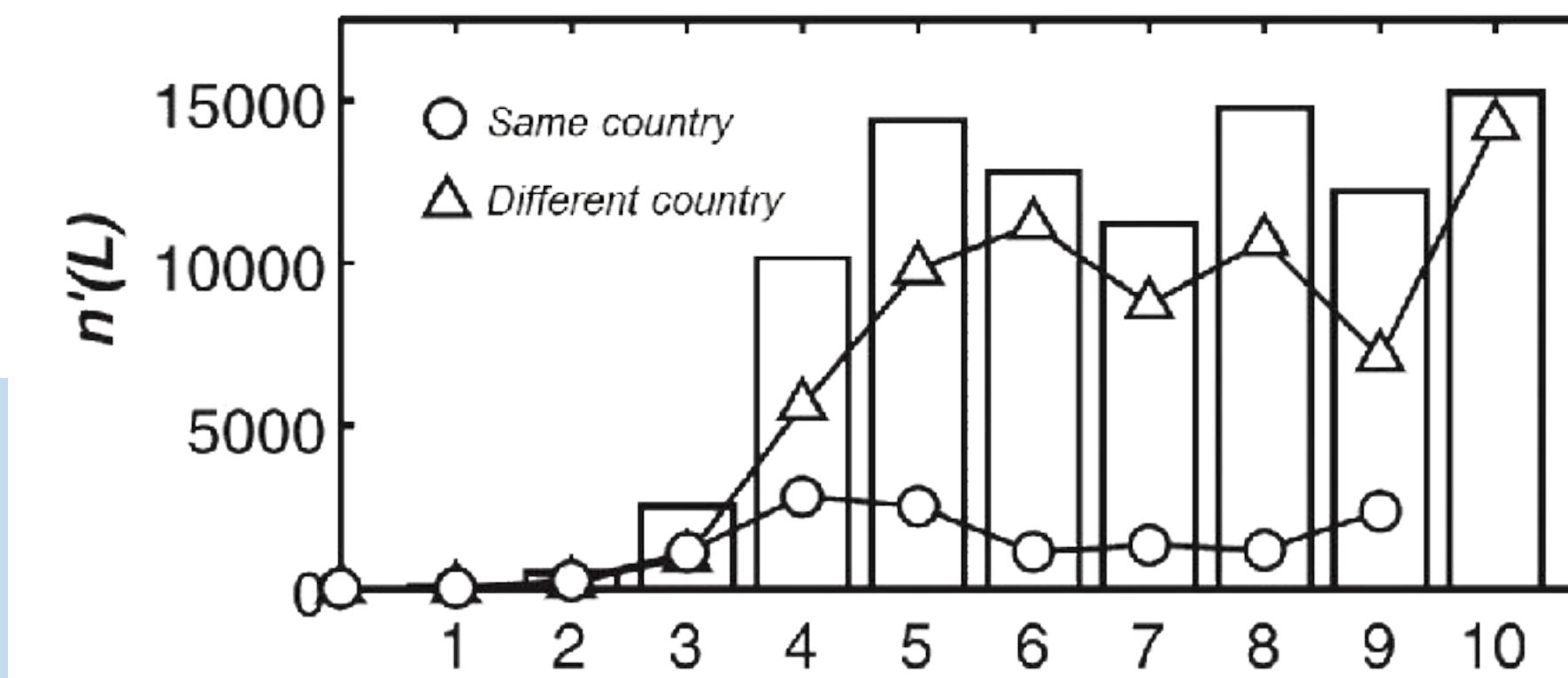


- Estimating length in a no-attrition scenario

$$n'(L) = \frac{n(L)}{\prod_{i=0}^{L-1} (1 - r_i)}$$

- Median number of steps = 7

“if individuals searching for targets do not have sufficient incentives to proceed, the small-world hypothesis **does not hold**”



Algorithmic vs. Topological definition

- Milgram explored the **algorithmic definition** of the small-world problem
 - Social search: shortest path that ordinary people can find given their **local topological information** about the underlying social graph
- **Large-scale web social data allowed for an exploration of the topological definition** of the same problem
 - Shortest possible path between two individuals, with a **global knowledge** of the social network
- See also “How small is the world, really?” by D. J. Watts
 - bit.ly/1oRZvIX

Topological small-world (MSN)

- **Availability of large-scale online social network data made direct topological measurement possible**
- MSN communication graph with 180M nodes and 1.3B undirected edges
- **Topological distance between 1000 random nodes and all others.**
Measure number of hops
 - Mode = 6
 - Median = 7
 - Average = 6.6
- https://arxiv.org/PS_cache/arxiv/pdf/0803/0803.0939v1.pdf

Topological small-world (Facebook)

- **Facebook data 721M nodes 69B edges**
- Compute distance between **all the pairs**
- Need a very efficient way to do it (HyperANF): <http://webgraph.dsi.unimi.it>
 - **Average of 4.74 hops**
- Recently repeated on **1.6B nodes**
 - **Average of 4.57 hops**
 - <https://research.fb.com/three-and-a-half-degrees-of-separation>

Open questions

- **Reliability**
 - Are all online social network ties good proxies for acquaintance?
- **External validity**
 - Does the result generalize across platforms?
 - Is the average path length shrinking in time?
- **Construct validity**
 - Is this results in line with expectations? What should we compare it to?
 - Are we measuring “how close the human population” is?
- **Ecological validity**
 - Is the actual graph distance meaningful for any social process?

Small Worlds

- We can characterize a network through 2 statistics:
- **Average shortest path length L**
 - this is not the same as the diameter of the graph, which is the maximum shortest path connecting any two nodes
- **The clustering coefficient C**
- A **small world graph** is any graph with a **relatively small L** and a **relatively large C**
 1. $L \propto \log N \approx L_{random}$
 2. $C \gg C_{random}$

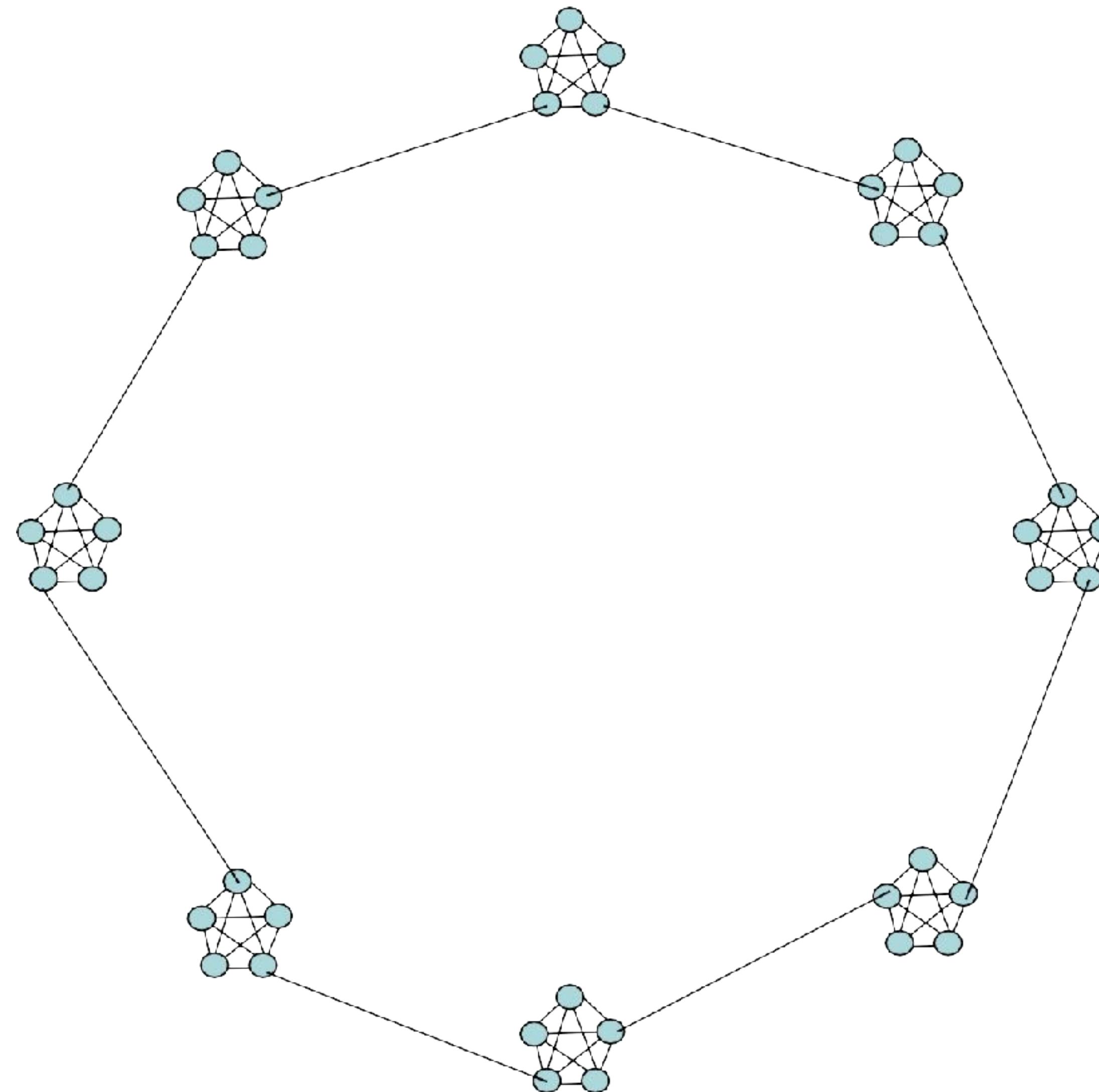
Found in many real-world phenomena

websites with navigation menus, food webs, electric power grids, metabolite processing networks, networks of brain neurons, voter networks, telephone call graphs, and social influence networks, cultural networks and word co-occurrence

Example

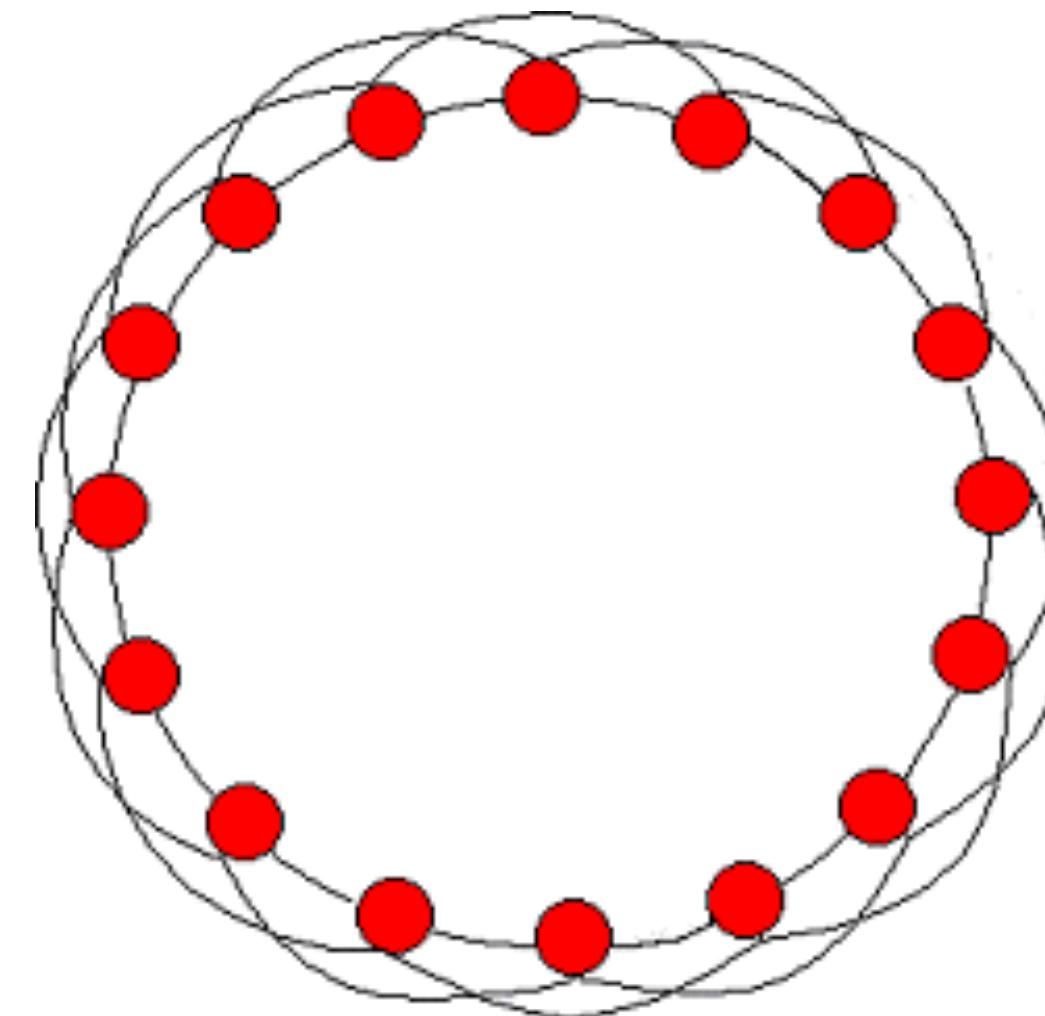
- Everyone in a cave knows each other
- A few people make connections
- **Are C and L high or low?**

C high
L high

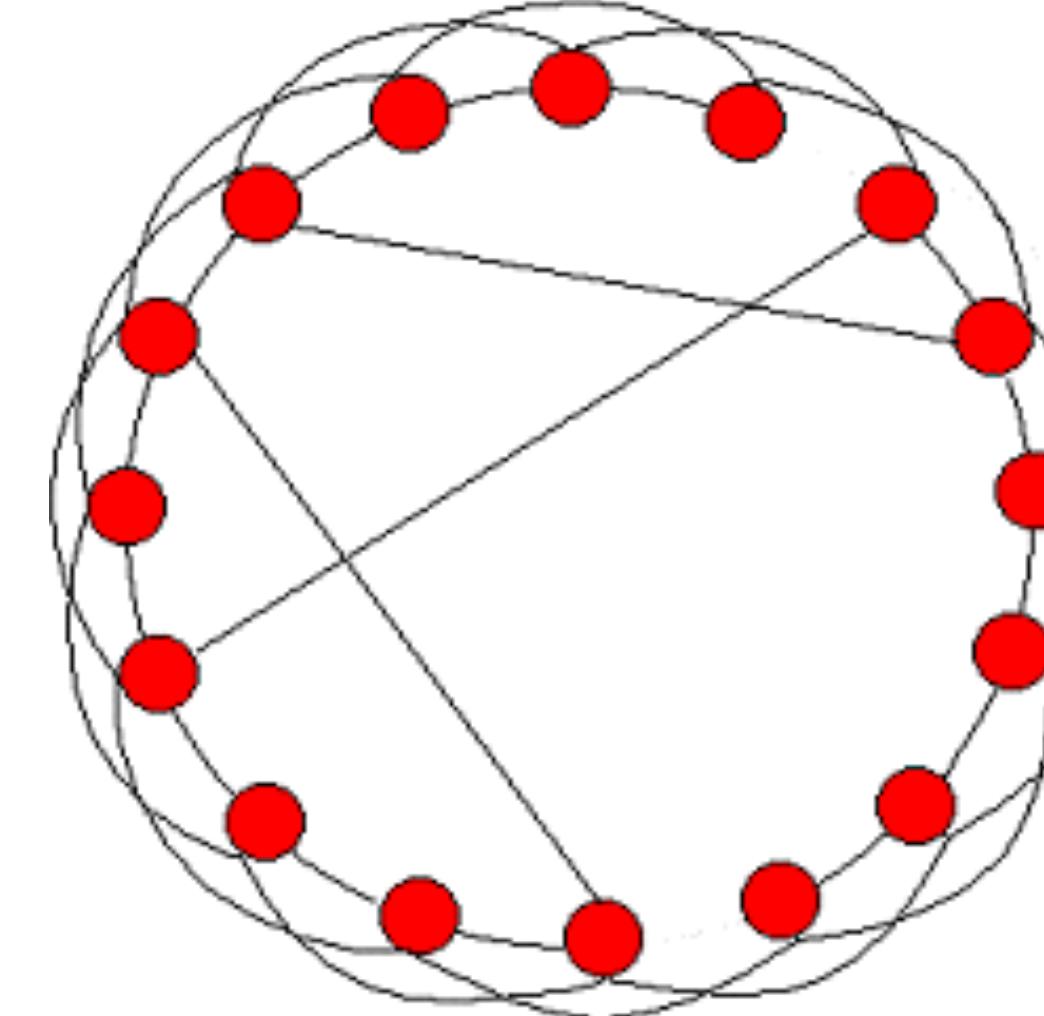


Watts and Strogatz model

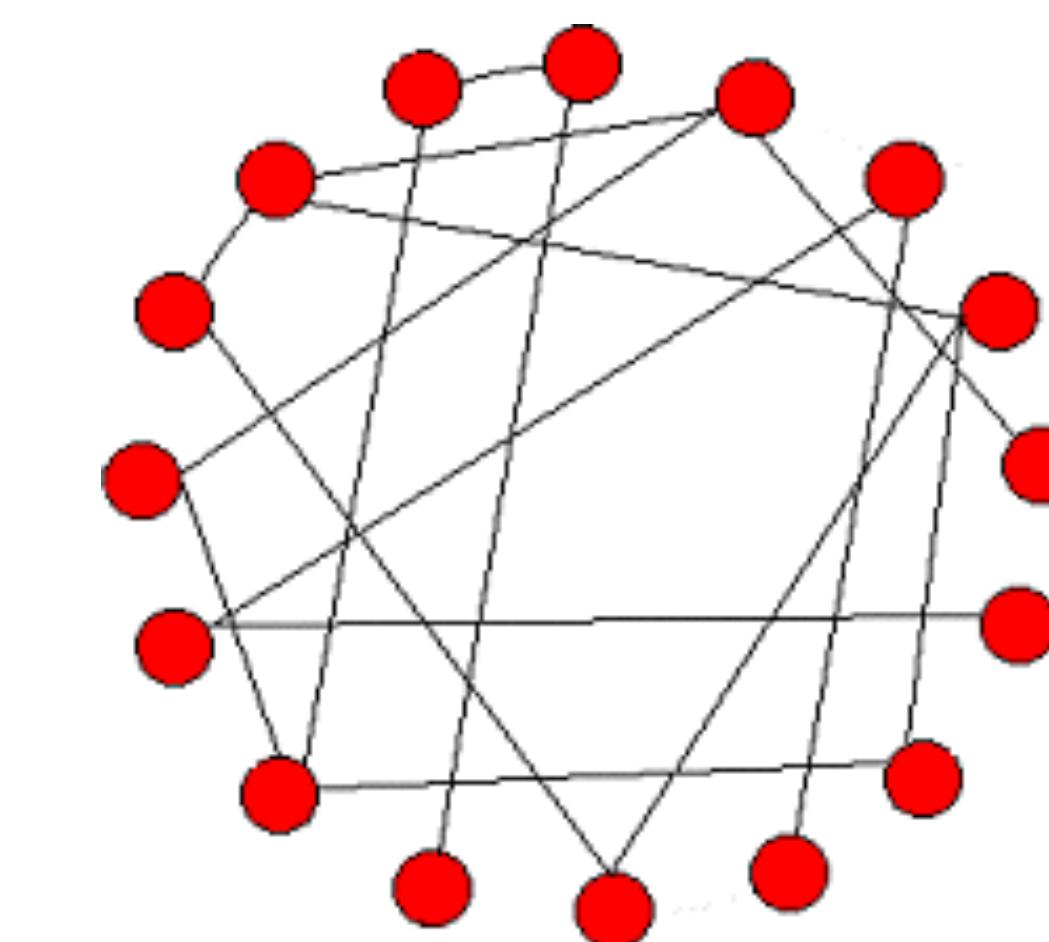
- Start with a ring, where **every node is connected to the next z nodes (a regular lattice)**
- **With probability α , rewire every edge (or, add a shortcut) to a uniformly chosen destination**



$\alpha = 0$
order



$0 < \alpha < 1$
Small world



$\alpha = 1$
randomness



Three Types of Networks

In a highly clustered, ordered network, **a single random connection will create a shortcut that lowers L dramatically**

Lattice

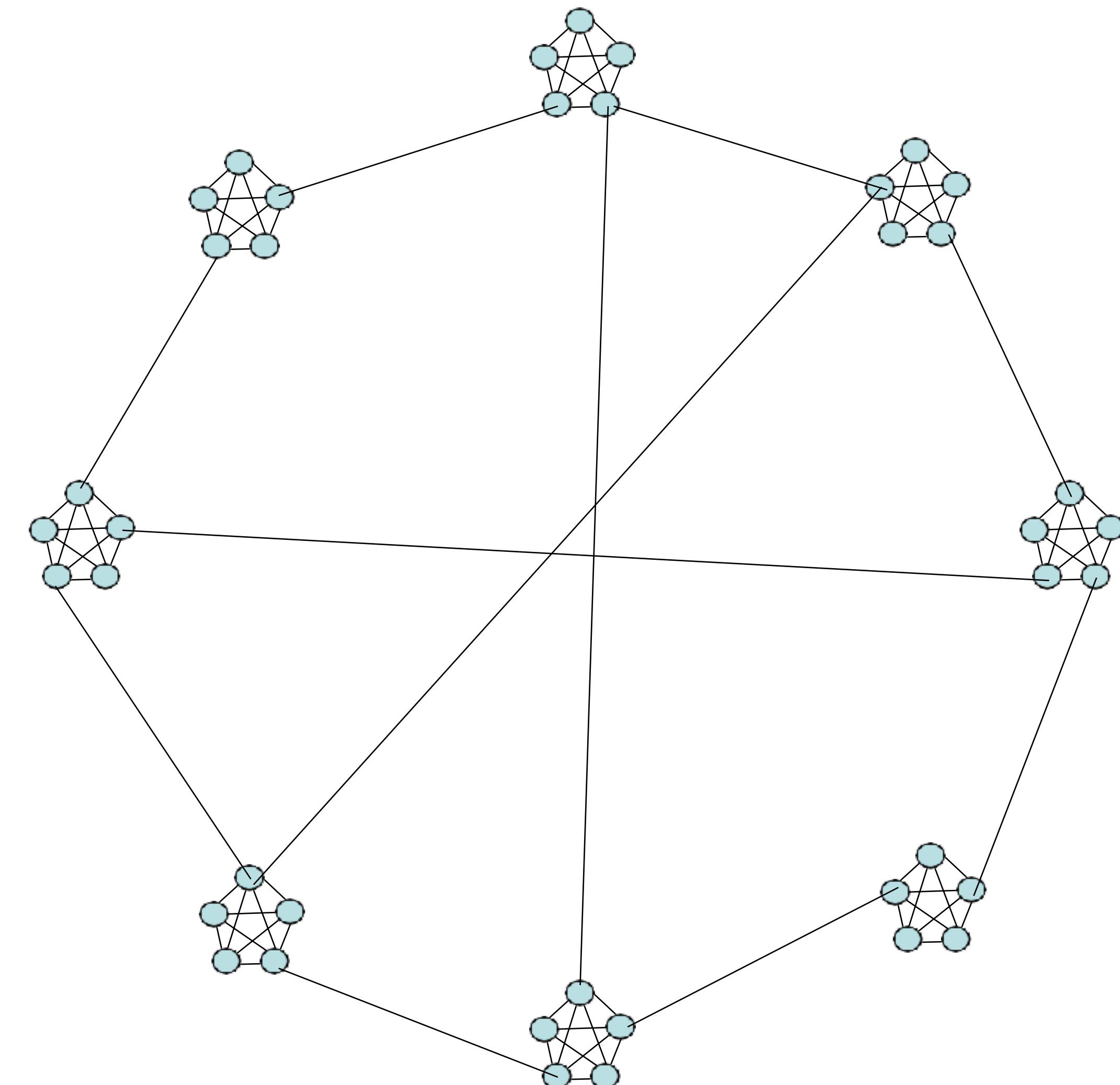
- spatial network where each person knows its neighbours
- **High clustering C**
- **High path length L**

Random Network

- neighbours of each node are randomly selected out of all nodes
- **Low clustering C**
- **Low path length L**

Small World Network

- **High clustering C**
- **Low path length L**



Tie Strength

Tie strength and social search

- In Milgram's Small-World experiment, participants had to report if the next step in the chain was a “**friend**” or “**acquaintance**”
- **Chain completion rate was 26% more likely if the first interracial tie was an “acquaintance”** (target person was black, start persons white)
- Sparked the intuition that **weak ties are more effective to exchange information over a social network**

Granovetter's tie strength

Paper was first rejected.
Now has 40k citations.

- The **degree of overlap** of two individuals' friendship networks varies directly with the strength of their social tie. **The strength of a tie is proportional to the similarity of its endpoints.**

“The **strength of a tie** is a (probably linear) combination of the amount of **TIME**, the emotional **INTENSITY**, the **INTIMACY** (mutual confiding), and the reciprocal **SERVICES** which characterize the tie.”

— Granovetter

Subjective and hardly measurable concepts!
But can be operationalized ...



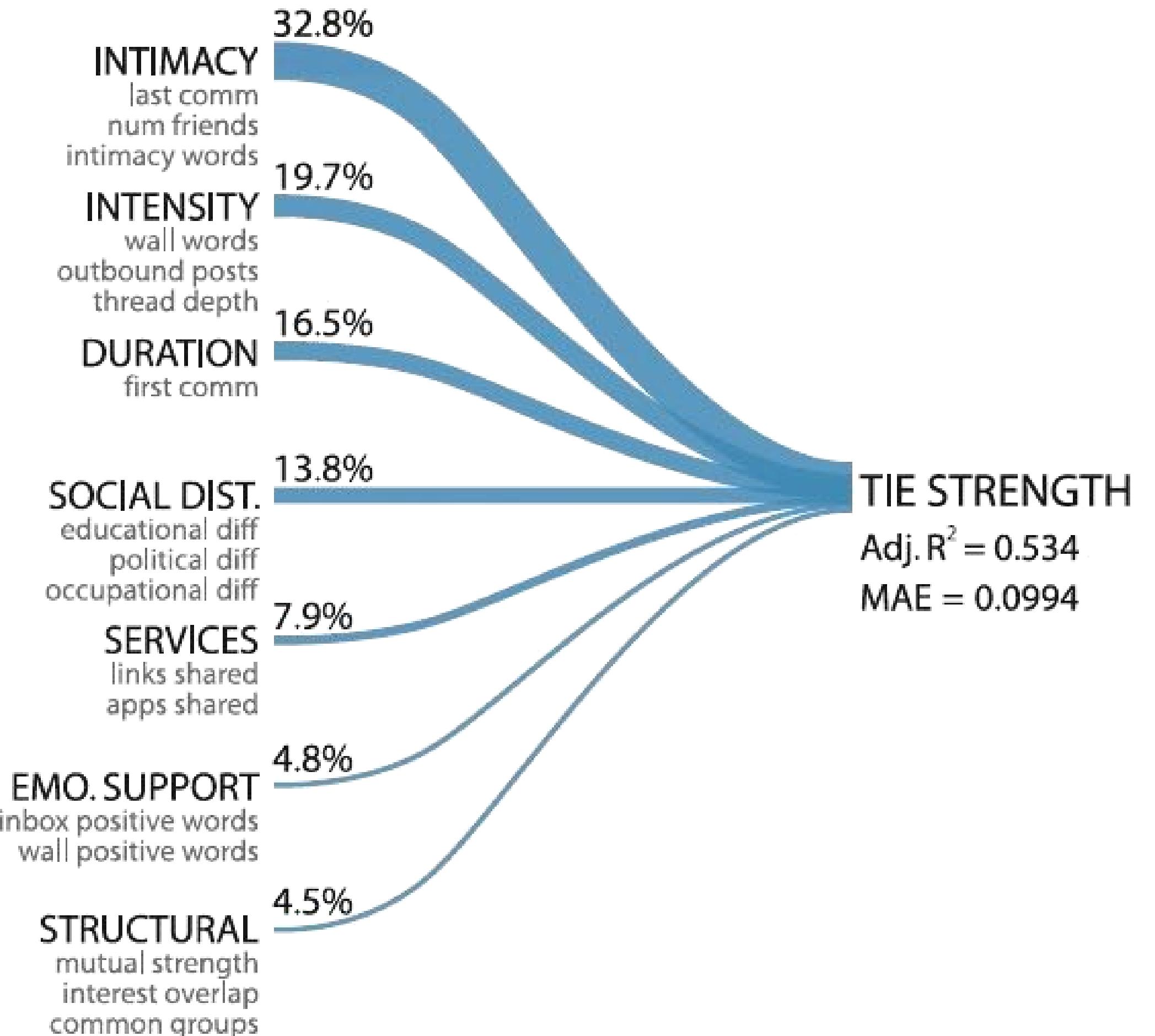
Definition



Operationalization

Objective measure of tie strength

- 2,184 self-rated Facebook friendship tie strength
- Goal: **predict tie strength** using proxies of Granovetter's indicators
- **85% prediction accuracy**



Granovetter's tie strength

- The stronger the tie between A and B
 - The higher the number of their shared social contacts
 - The more similar they are

A+B and A+C are friends
A+B spend time together
A+C spend time together

→ C+B will be acquainted

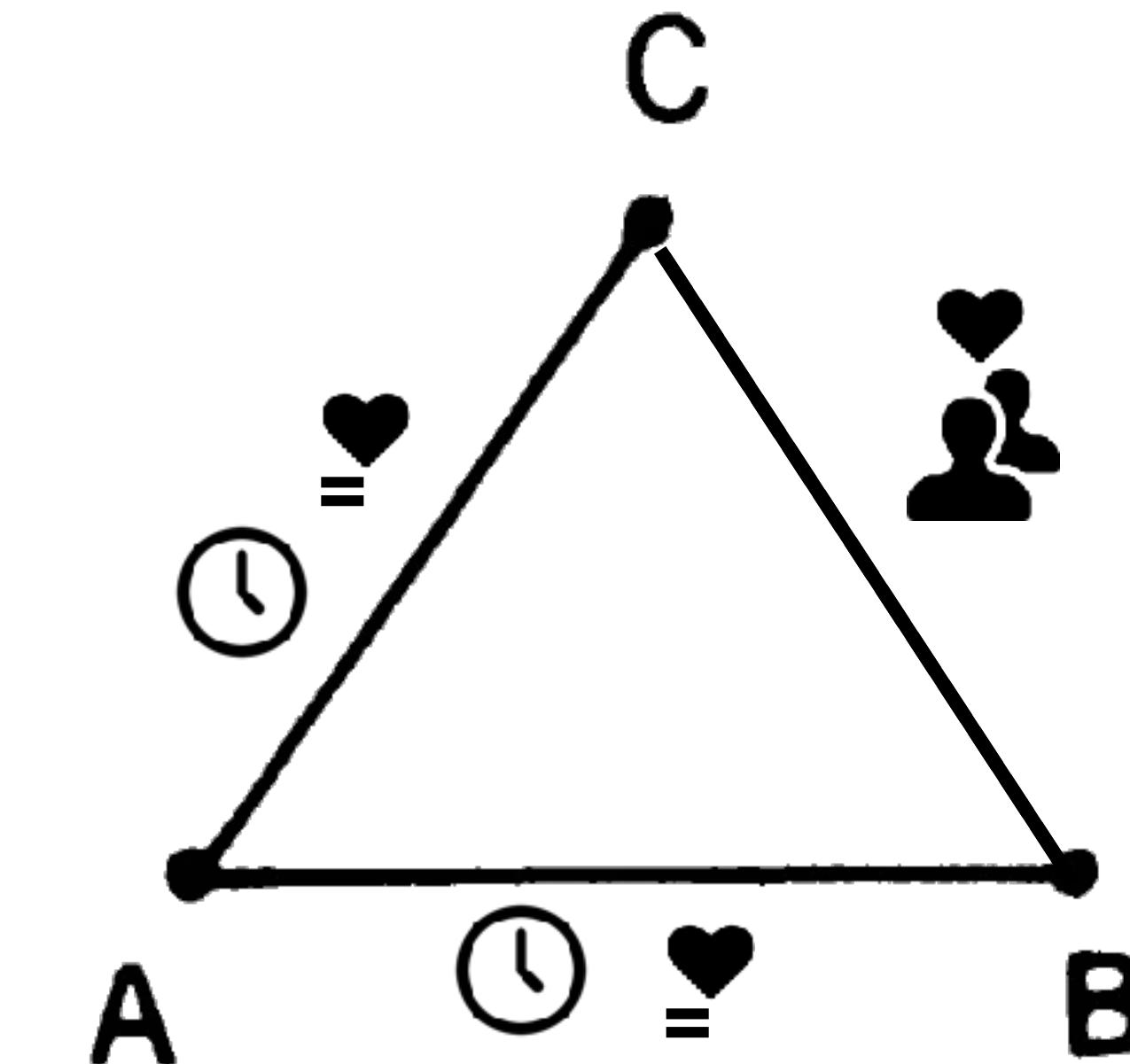
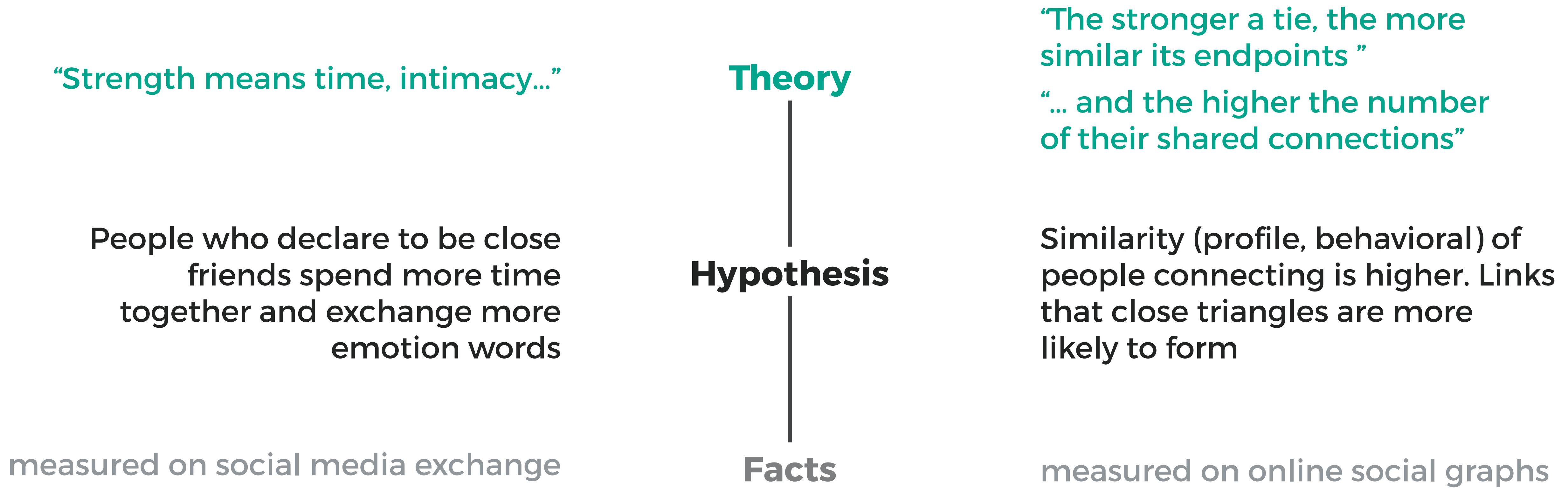


FIG. 1.—Forbidden triad

Time spent = tie strength

Theory -> Hypothesis -> Facts

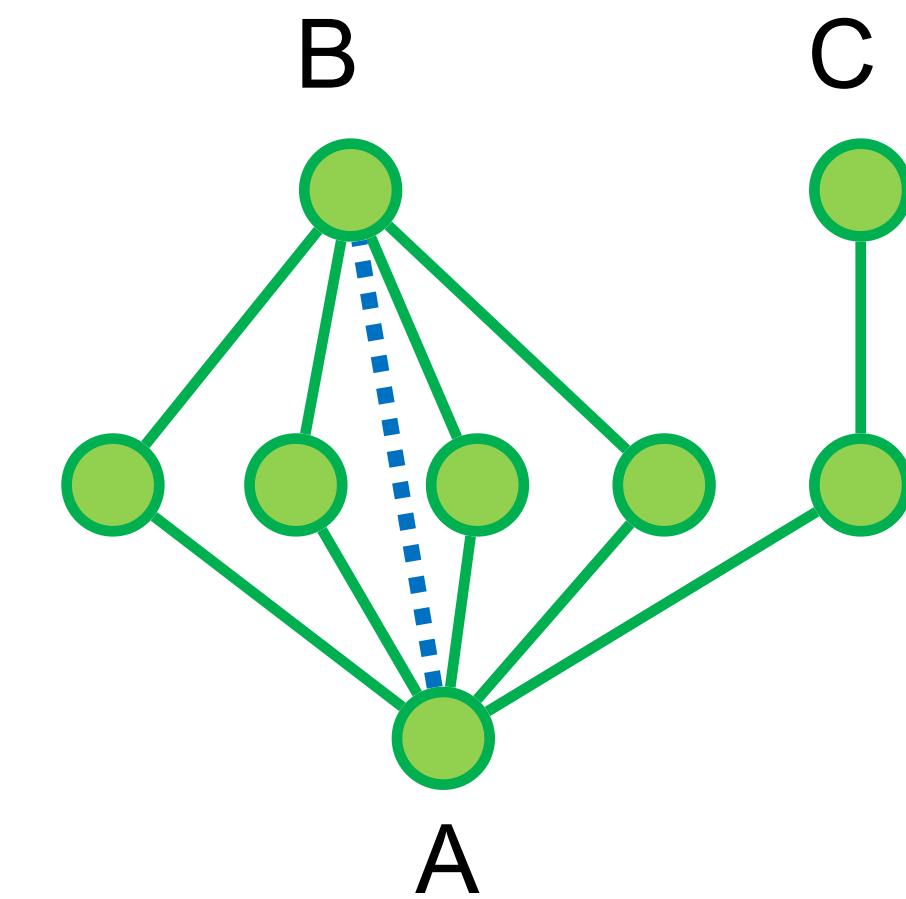


Validity? Do more (and different) experiments

Hypothesis 1: Triangle closure

structural

- **The more triangles the edge closes, the higher the likelihood of its formation**
- Can be tested with **predictive frameworks**
 - Target node A
 - Its distance-2 neighbors at time t: $\Gamma^2_t(A)$
 - Predict the nodes in $\Gamma^2_t(A)$ that will be connected with A at time $t+\Delta$
- **Number of common neighbors (and derivative measures) is highly predictive**
 - e.g., Jaccard coefficient



Hypothesis 2: Similarity

homophily



Activity patterns of users are highly heterogeneous

k-in in the friendship network

k-out in the friendship network

n_t of distinct tags in the vocabulary

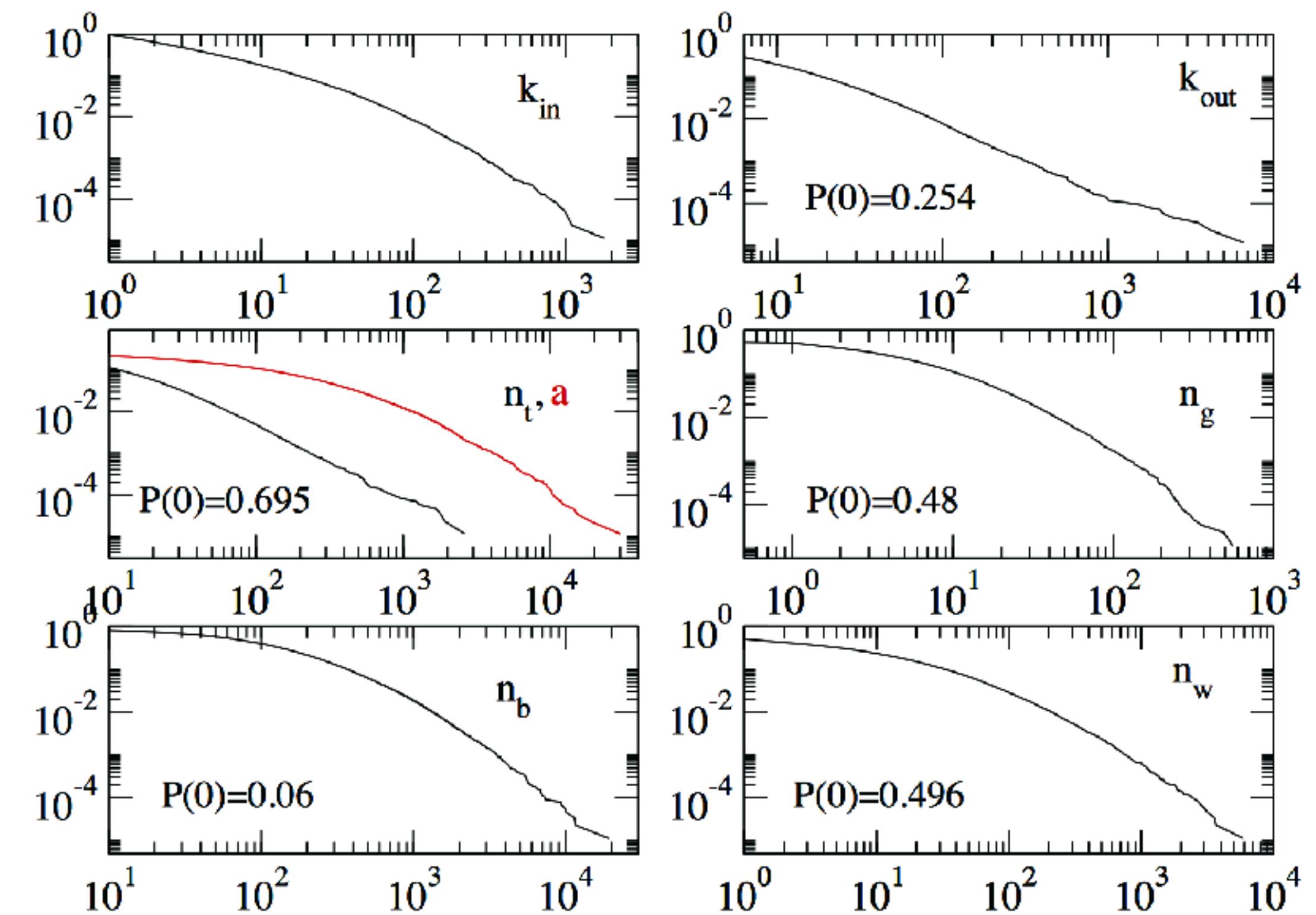
n_g groups

a total tagging activity

n_b number of books in the library

n_w number of books in the wishlist

Activity metrics are all positively correlated

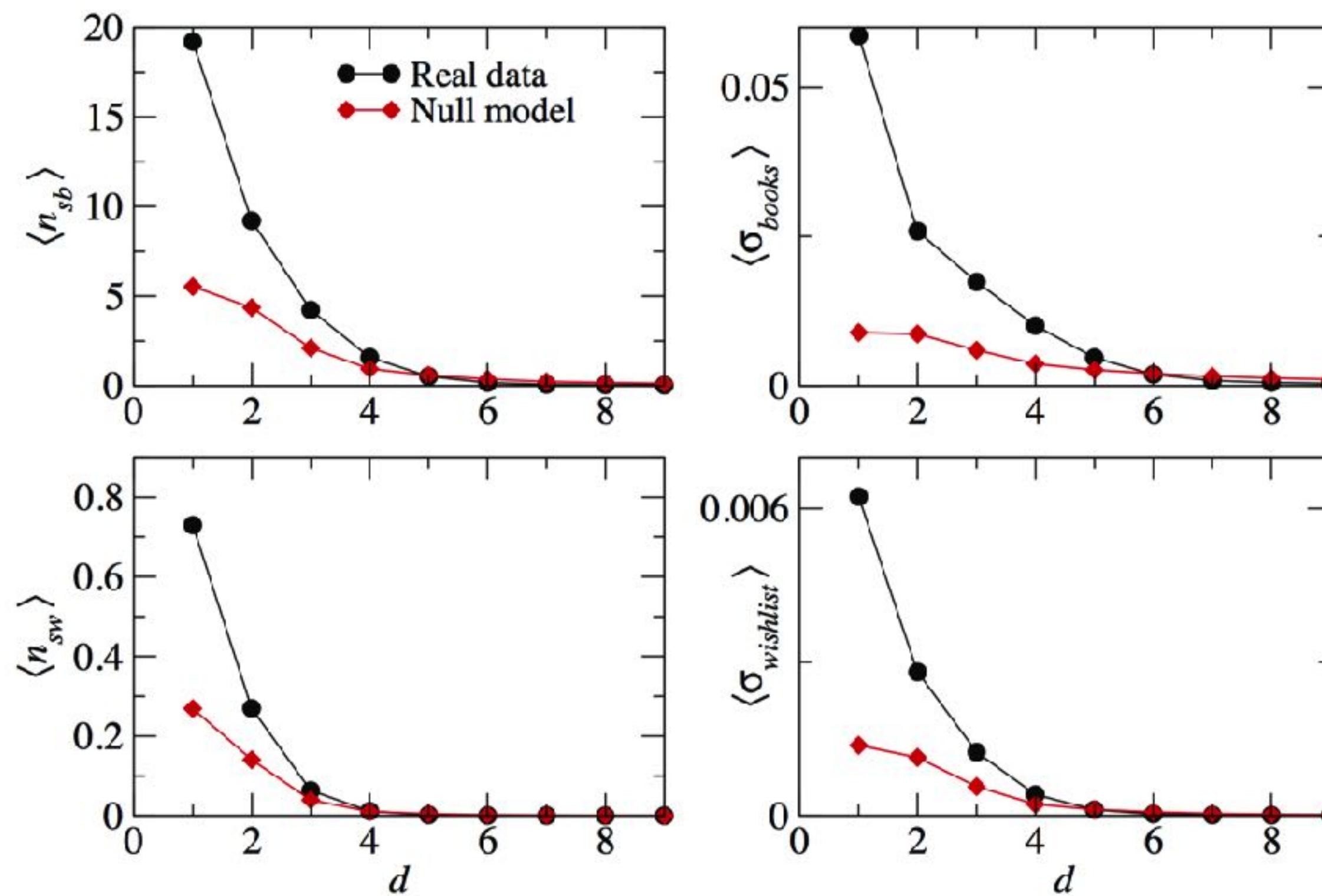


Hypothesis 2: Similarity

homophily



Topical Similarity



n_{sb} number of shared books in the library

n_{sw} number of shared books in the wishlist

δ_{books} cosine similarity books vector in the library

δ_{wishlist} cosine similarity books vector in the wishlist

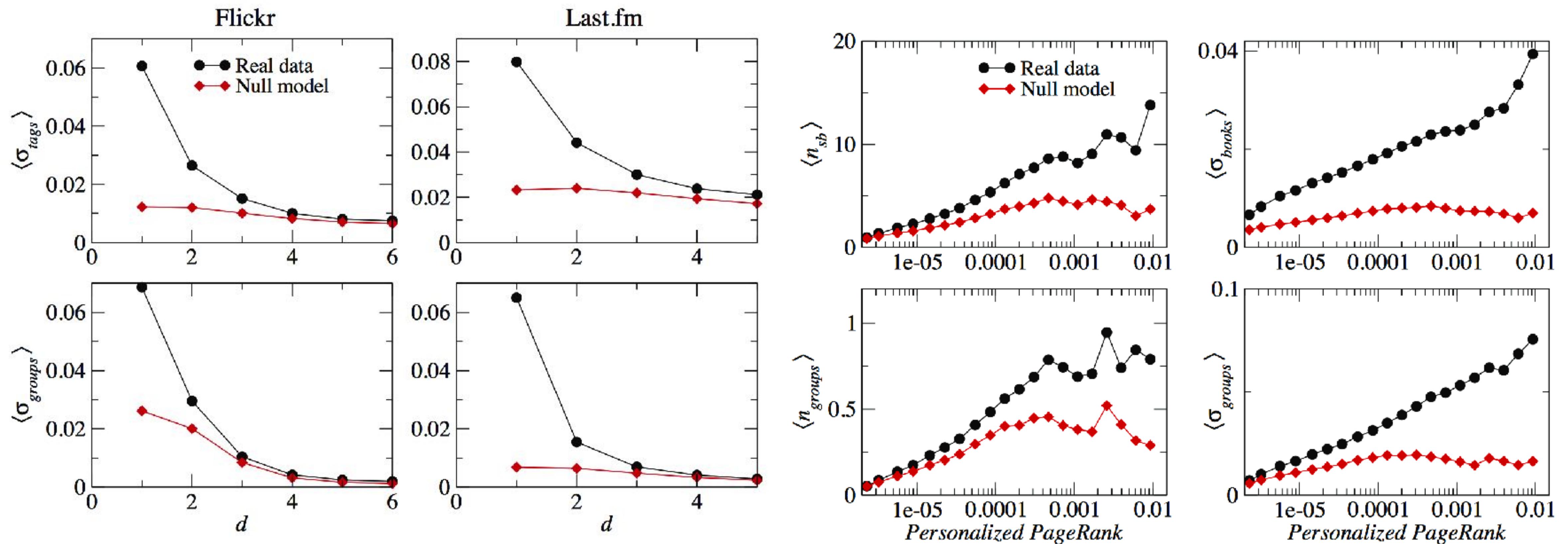
δ_{tags} cosine similarity tags vector

δ_{groups} cosine similarity groups vector

Hypothesis 2: Similarity

homophily

Topical Similarity



Link creation and information spreading over social and communication ties in aNobii, EPJ Data Science 2012
Friendship prediction and homophily in social media, TWEB 2012

Hypothesis 2: Similarity

Link prediction

homophily



Feature	Similarity	Last.fm		aNobii	
		Active	Connected	Active	Connected
Baselines	Tasteometer	0.734	0.759	-	-
	Common neighbors	0.927	-	0.854	-
Items	Distrib cosine	0.663	0.560	0.915	0.655
	Distrib MIP	0.749	0.559	0.878	0.649
	Collab MIP	0.589	0.613	0.652	0.561
Tags	Distrib cosine	0.579	0.625	0.652	0.554
	Distrib MIP	0.697	0.618	0.651	0.560
	Collab MIP	0.698	0.559	0.916	0.648
Groups	Cosine	0.810	0.677	0.662	0.690
Library	Cosine	-	0.769	0.923	0.768

Takeaway

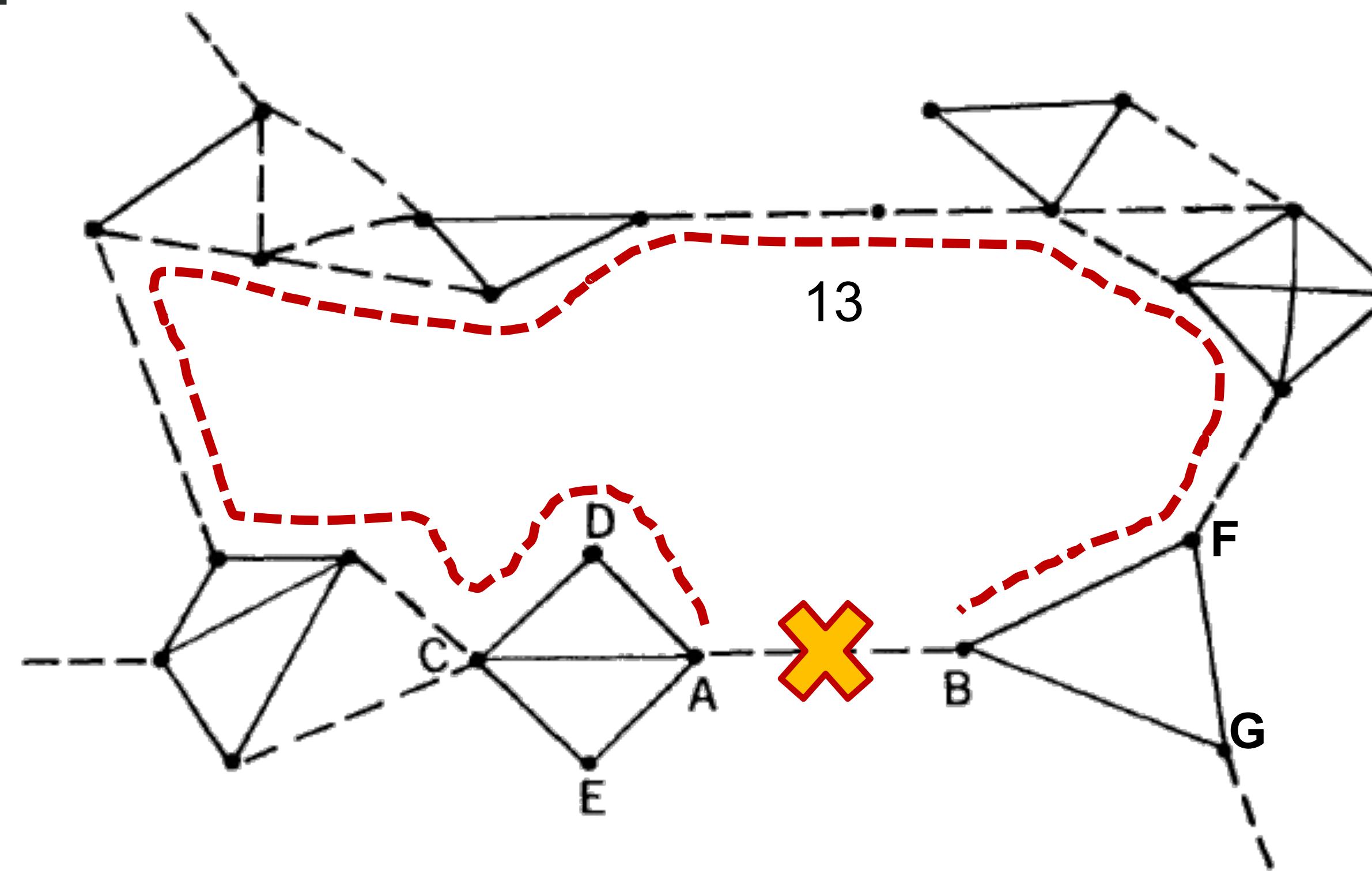
- Similarity on profile and activity is a good predictor**
- Structural + Profile features perform the best**

	Tags	collab	MIP	Groups	Library	CN	Profile feaures	All features
AUC		0.785		0.807	0.811	0.844	0.924	0.963
Accuracy		0.786		0.809	0.812	0.846	0.877	0.915
FP Rate		0.177		0.143	0.335	0.031	0.109	0.072
FN Rate		0.251		0.240	0.041	0.277	0.137	0.099

Link creation and information spreading over social and communication ties in aNobii, EPJ Data Science 2012
Friendship prediction and homophily in social media, TWEB 2012

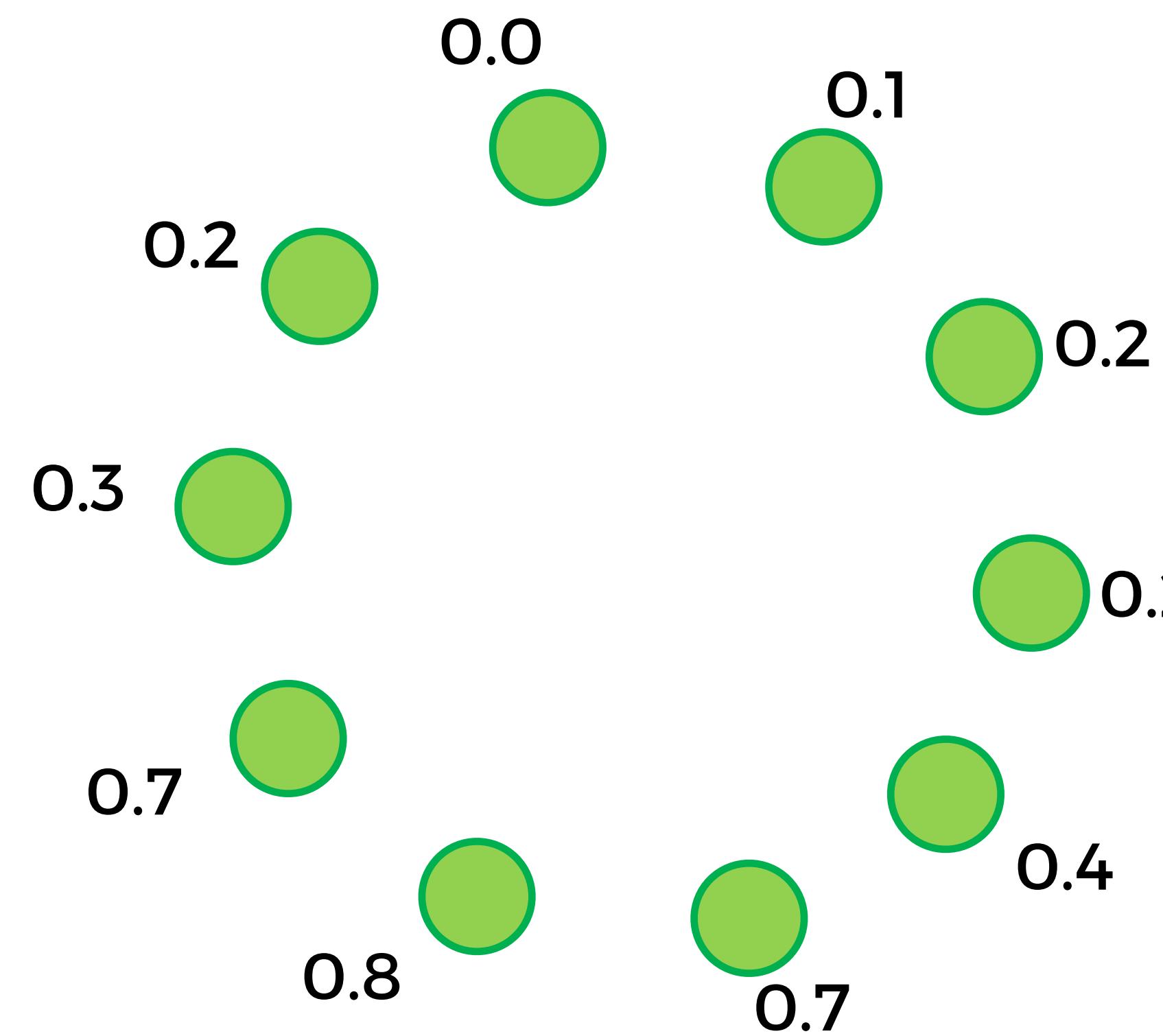
Corollary: Weak ties and bridges

- **Bridges are always weak ties**, never strong ties
- **Strong ties bring fragmentation**: information has way more pathways to stay within dense local clusters
- **Weak ties reduce path length** and create **opportunities to access new information**



Threshold model of diffusion

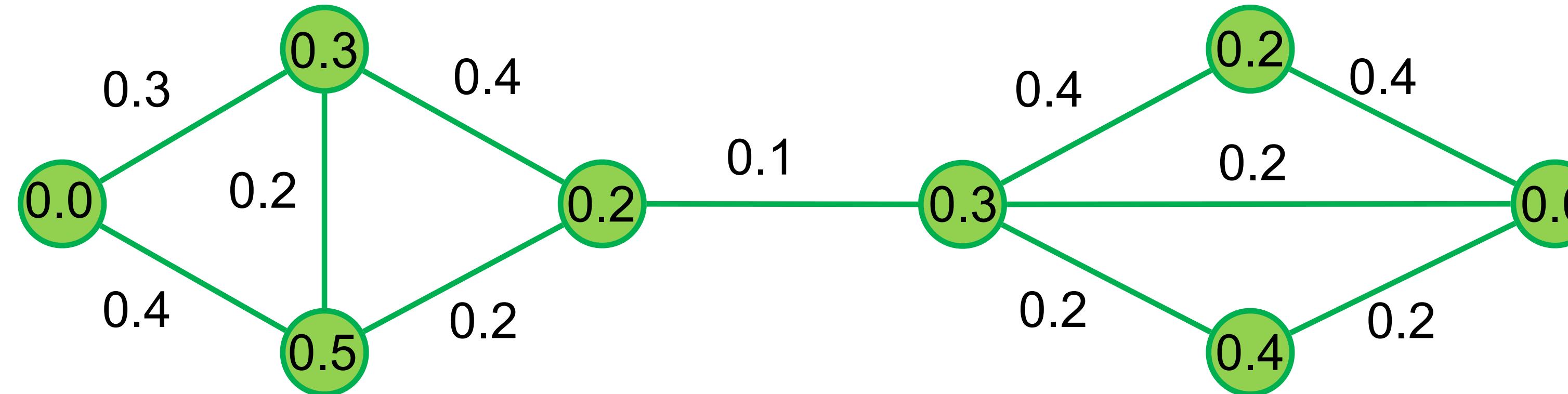
- Individuals with binary activation status
- Each person i has an activation threshold $\Phi(i)$ in $[0,1]$
- At time t , i activates if the fraction of active nodes is $\geq \Phi(i)$



Complex contagion (multiple exposures) vs. simple contagion (contact with disease)

Linear model threshold

- Tie strength and individual threshold can be combined in a networked model
- Each person i has an **activation threshold** $\Phi(i)$ in $[0, 1]$
- At time t , i activates if sum of strengths of ties towards active people is $\geq \Phi(i)$



Another example

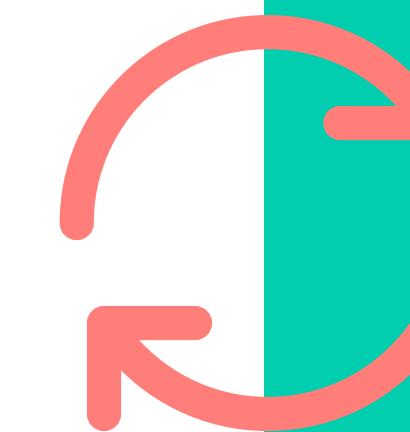
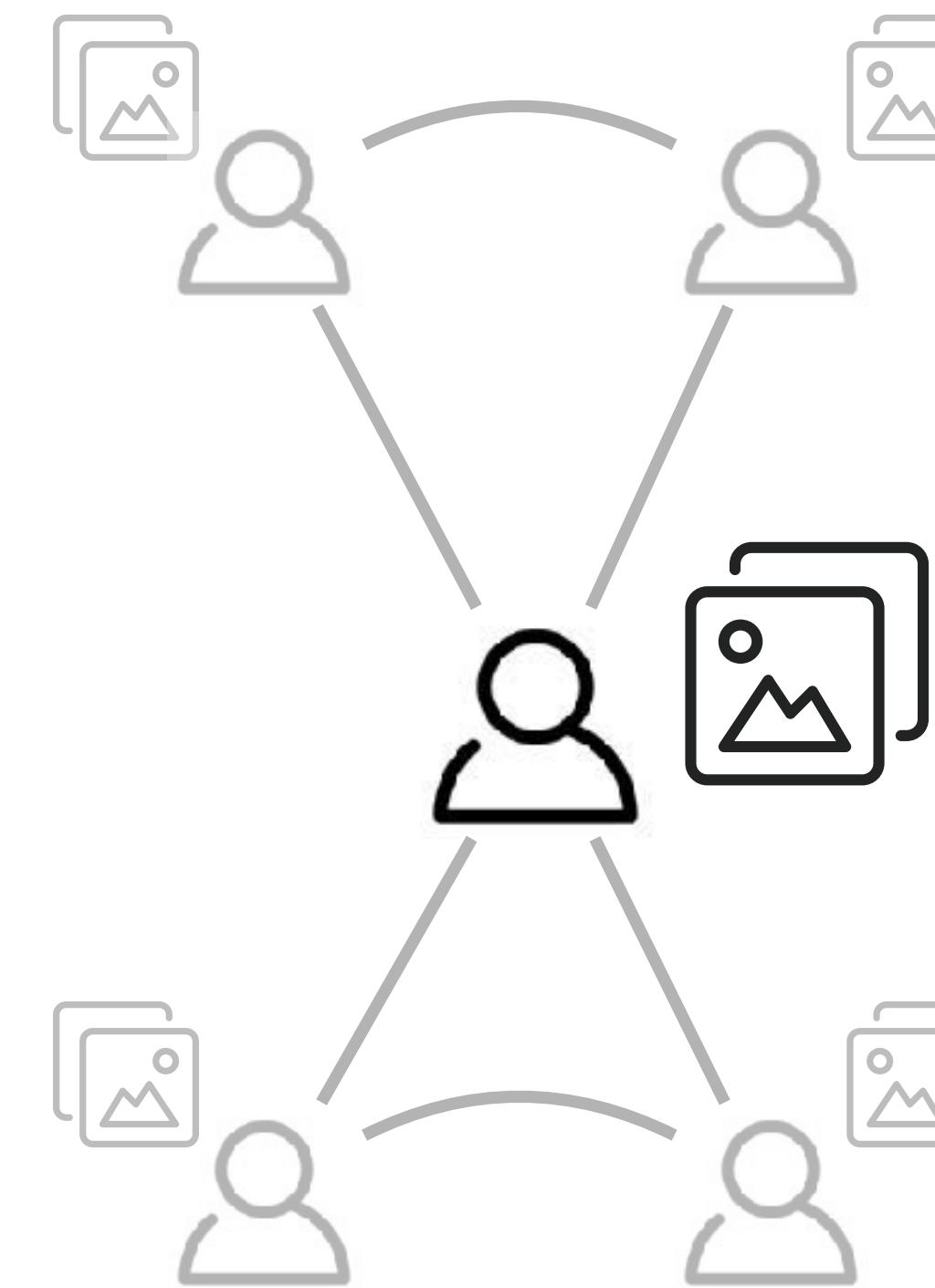
- In a photo sharing social network (Instagram), **people post pictures of the icecreams** they are eating. When users upload a photo, all their friends are notified. Friendship links are directed according to the information flow: a directed edge A->B means that B follows A and that B is notified with A's profile updates.
- We observe a series of "icecream consumptions" (in the form [userA, timeX]) over a period of time. Some users might consume icecream multiple times. The social network is static over that period of time. **We would like to check whether the icecream consumption is driven by a phenomenon of "social contagion".**
- **Hypothesis: People consume an icecream when more than x% of their friends have recently consumed an icecream (activation threshold).**

Social Influence, Quality and Engagement in Social Media

social
popularity
(success)

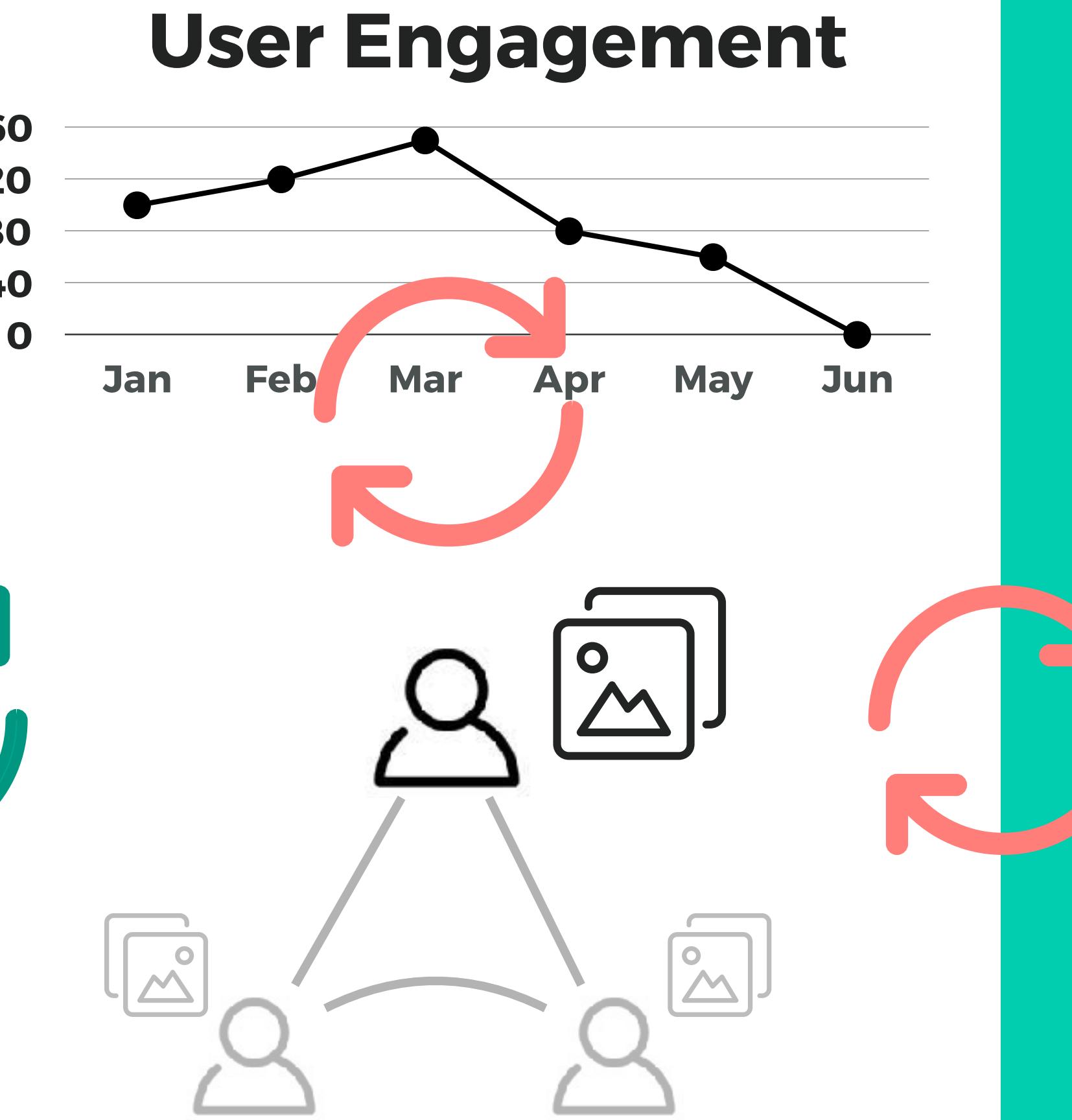


social network



intrinsic
quality
(performance)

social
popularity
(success)



intrinsic
quality
(performance)

social network

Questions

1. Popularity = Quality ?
2. Do social ties affect quality?
3. Are network effects intertwined with engagement?

social popularity (success)

YouTube ACTIVITY (LIKES, VIEWS, COMMENTS)



CONNECTIVITY (FOLLOWERS, FRIENDS)
WITNESS THE TOUR



KATY PERRY 

@katyperry
katyperry.com/tour
Joined February 2009
Born on October 25

Tweets 8,730 Following 204 Followers 106M Likes 5,817

[Follow](#)

Tweets [Tweets & replies](#) [Media](#)

Pinned Tweet
KATY PERRY  @katyperry - May 15
• WITNESS • THE ALBUM ! THE TOUR ! IT'S ALL HAPPENING !
katyperry.com/WITNESSKP

Who to follow [Refresh](#) [View all](#)

- Rihanna  @rihanna [Follow](#)
- Justin Timberlake  @jtim... [Follow](#)
- Lady Gaga  @ladygaga [Follow](#)

[Find people you know](#)

Trends for you [Change](#)

Visual Content



intrinsic
quality
(performance)

Visual Aesthetics

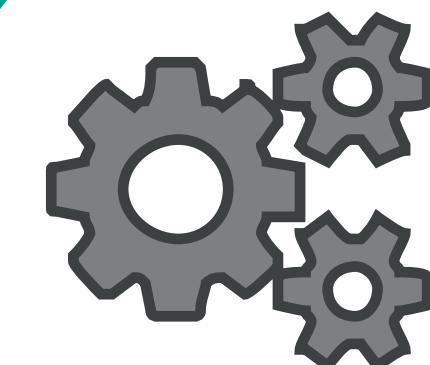
- **Multidisciplinary**
 - **Philosophy**
 - **Experimental Aesthetics (psychology)**
 - complex interaction between perception, cognition and emotion
 - **Neuroaesthetics (neurology)**
 - **Computational Aesthetics (computer vision)**

How beautiful is this photo?



human score

4.5



Computational
Aesthetics Engine

machine score

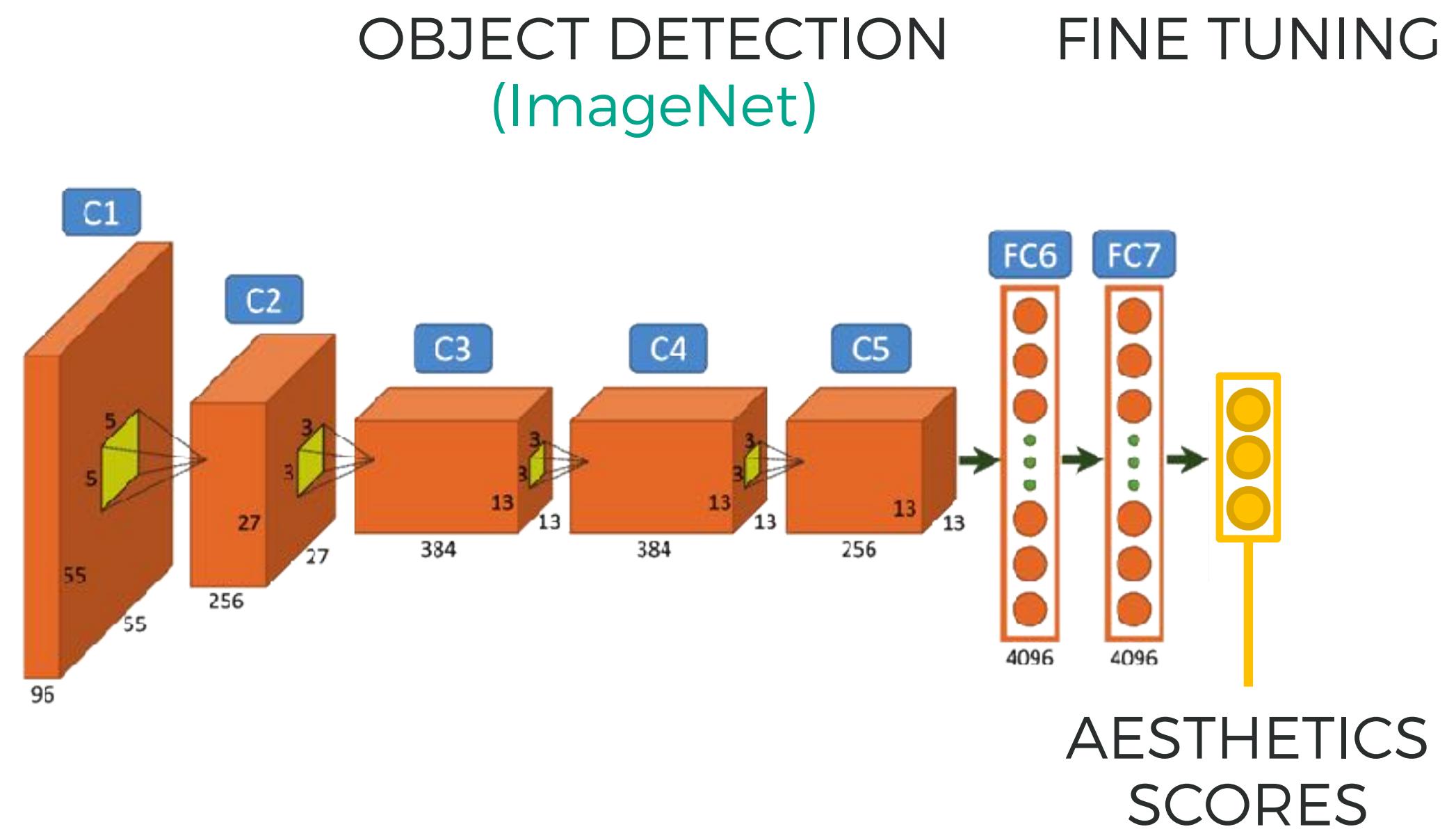
4.3

Computational Aesthetics

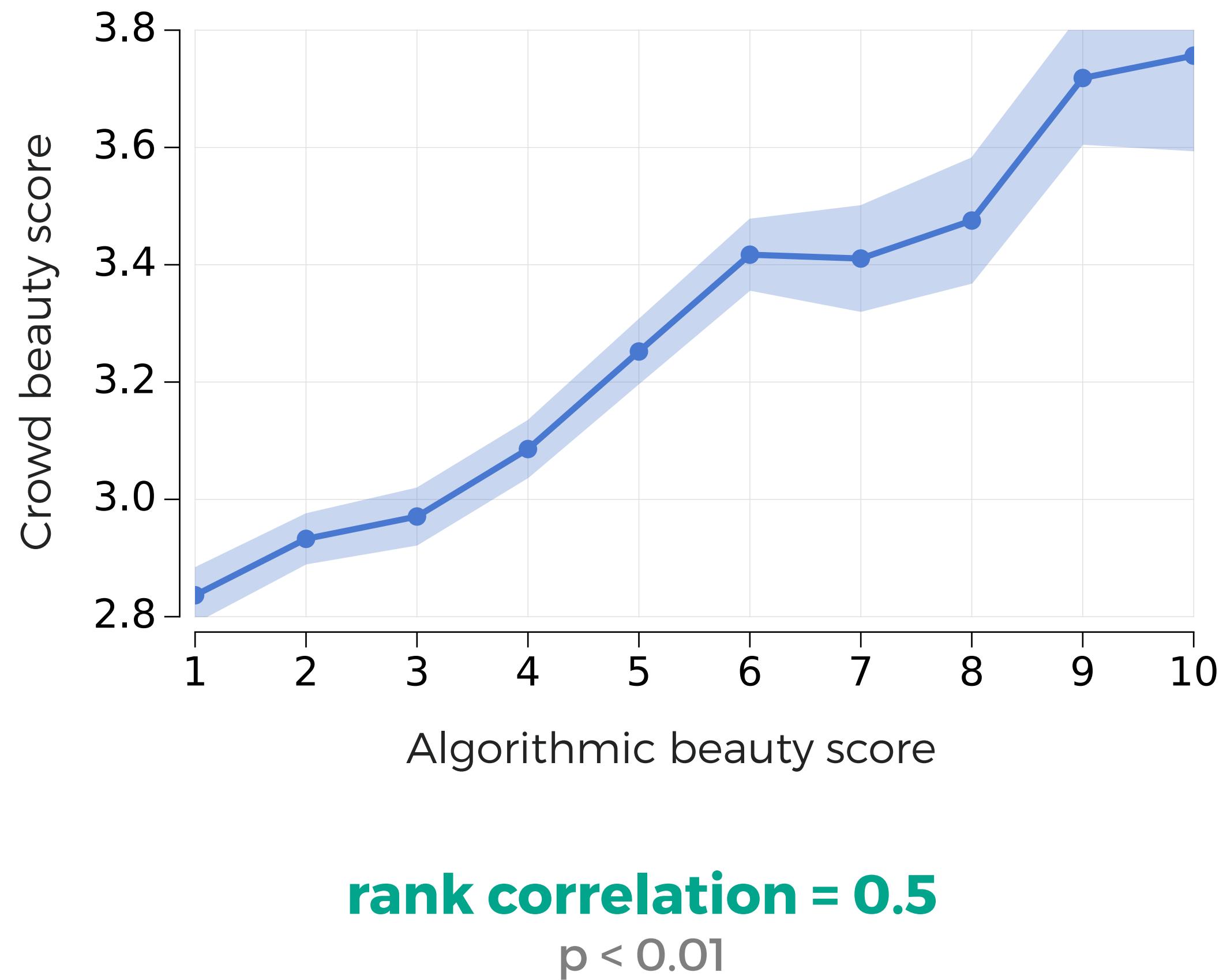
- **Hand-Crafted features**
 - **low level**
 - blur, contrast, colorfulness, saliency, texture, shape, energy
 - **compositional**
 - rule of thirds or depth of field
- **Generic Image Features**
 - **deep learning**
 - content, form, and [blackbox]

Algorithmic Beauty

(Transfer Learning)



Machine eye ≈ Human eye



Human evaluation

- **4 Categories**
 - nature, animals, people, urban
- **Full spectrum of popularity**
- **Good agreement**
- **Large ground truth of aesthetic scores made by non-expert**
 - 10K photos
 - 60K scores
- **Available at [http://di.unito.it/
beautycwsm15](http://di.unito.it/beautycwsm15)**

1.

POPULARITY AND QUALITY DO NOT (ALWAYS) ALIGN

aNobii.com

 **Aliceassassina**
Female, 34. Roma, Italy

[Follow](#) [Message](#)

Friends [more](#)

-  [La Loulu](#)
-  [Tripwood76](#)
-  [Nerorossobia...](#)

Neighbors [more](#)

-  [Zoro](#)
-  [maldido duen...](#)

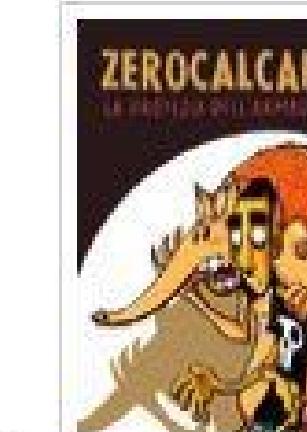
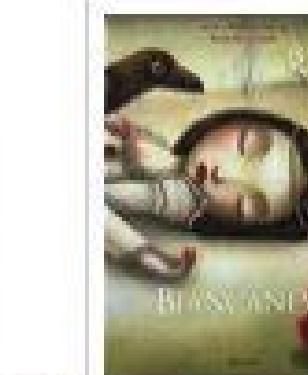
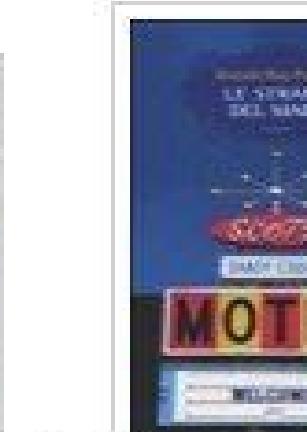
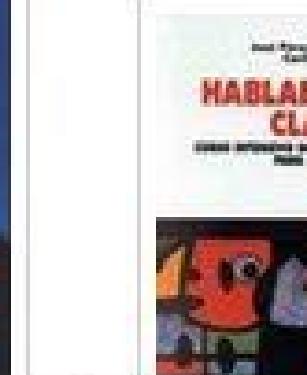
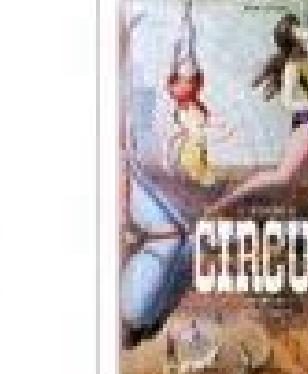
Shoutbox [See all](#)

Leave a comment

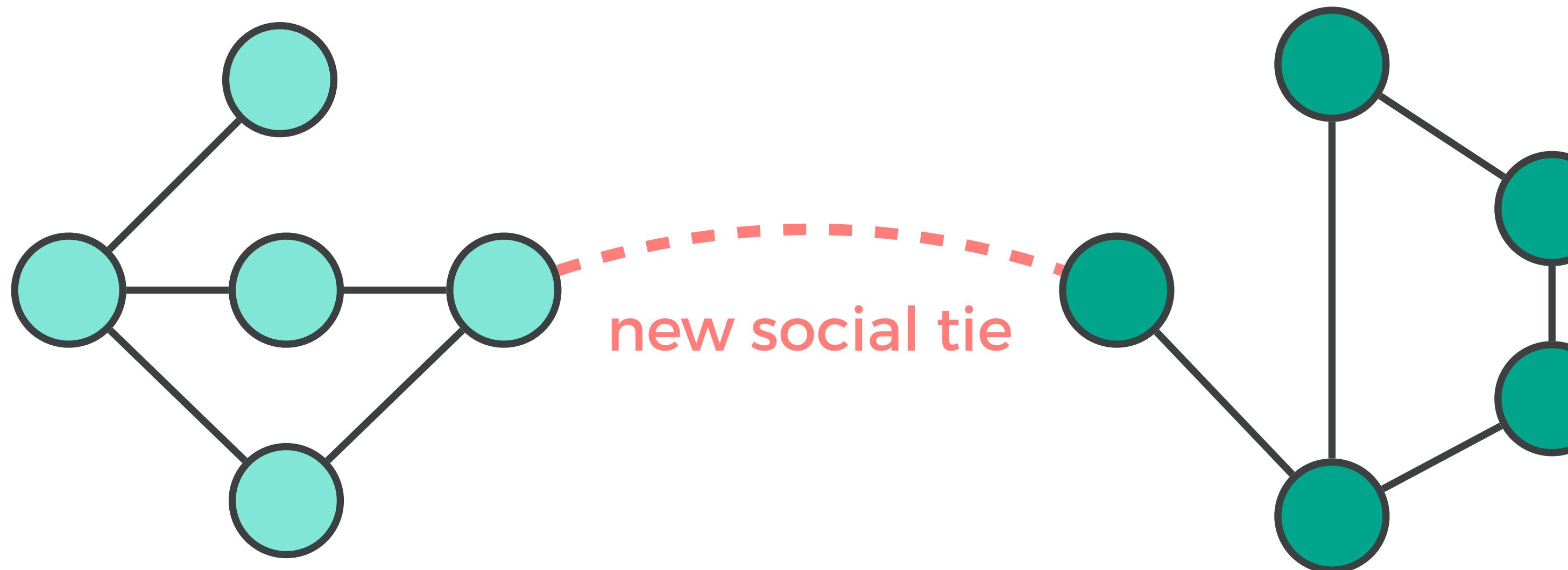
Morgana1981 "16 Novembre 2013 è uscito "Le Lenzuola nere di Khloe", Oct 28, 2013
CIAO!
Questo è il booktrailer del mio nuovo libro: "Le Lenzuola Nere di Khloe". Dura 2 minutini scarsi!
Buona Visione! :)

Books

All books ▾ Rating (Highest) Search this shelf

 L'Erbario delle Fate (85) By Benjamin Lacombe, Sébastien Perez Finished on Nov 28, 2013 ★★★★★ comment + Add...	 La profezia dell'armadillo (2931) By Zerocalcare Finished on Jan 6, 2013 ★★★★★ comment + Add...	 Biancaneve (171) By Jacob Grimm, Wilhelm Grimm Finished on Feb 11, 2012 ★★★★★ comment + Add...	 Le strade del male (107) By Donald Ray Pollock Finished on Aug 17, 2012 ★★★★★ comment + Add...	 Hablando claro (14) Curso intensivo de Espanol para italianos By Carla Polettini, José Pérez Navarro Finished in 1999 ★★★★★ comment + Add...
 La ballata di Trenchmouth (50)	 C'era una volta... (113) Libro pop-up	 Circus Book 1870 - 1950 (4) Reference	 Papà tatuato (67) By Daniel	 Venere privata (2507) By Giorgio

Our goal was to predict links

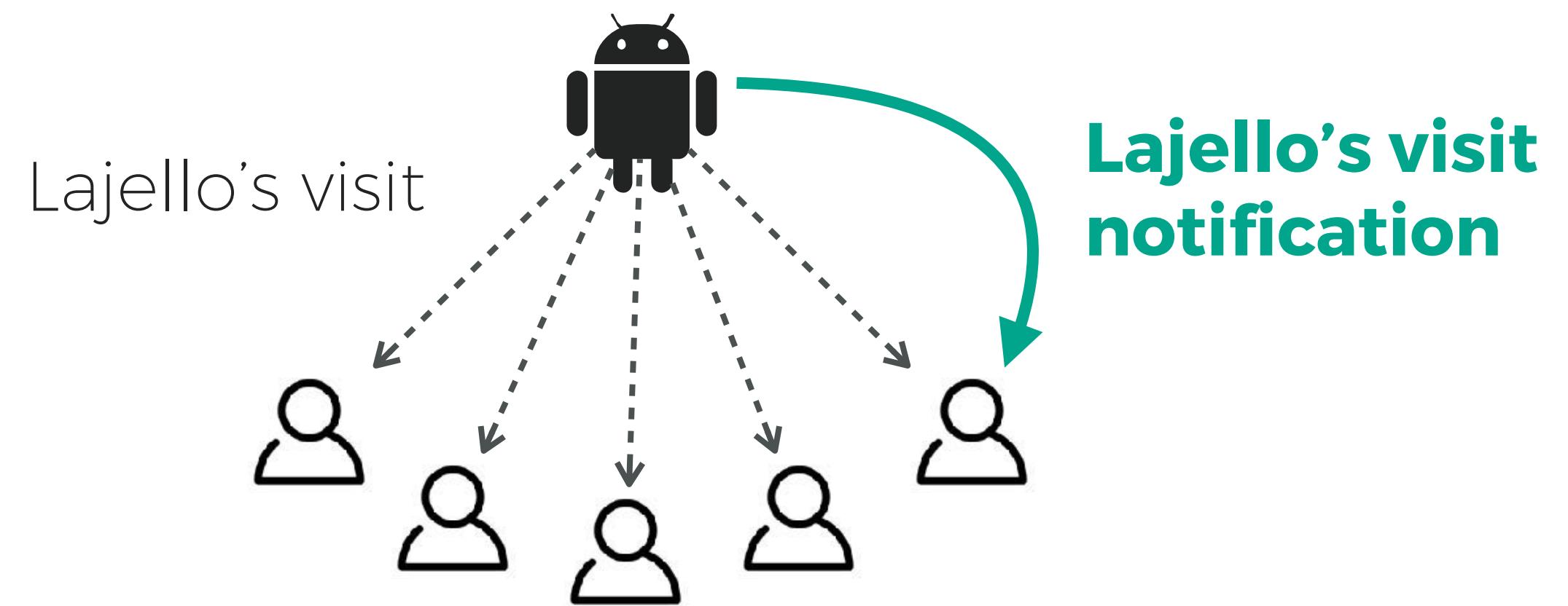


Lajello crawler



totally empty profile!

Iterative probing
(every week)



Gain of popularity

2ND
(OUT OF 100K)

>1.2K
DISTINCT USERS

28TH
(OUT OF 100K)

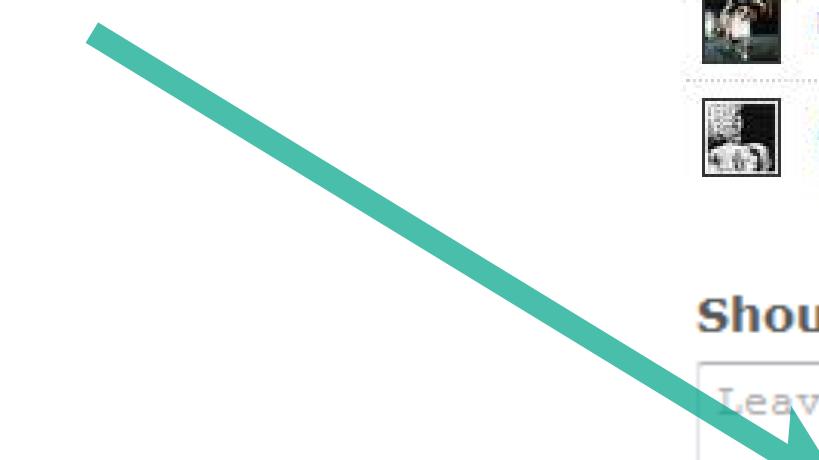
~2.5K
PUBLIC MESSAGES

0.3%

>200
INCOMING LINKS

Influence

Lajello left a personalized comment in the **shout box** of 2K users



 **Aliceassassina**
Female, 34. Roma, Italy

[Follow](#) [Message](#)

Friends [more](#)

-  [La Loulu](#)
-  [Tripwood76](#)
-  [Nerorossobia...](#)

Neighbors [more](#)

-  [Zoro](#)
-  [maldido duen...](#)
-  [ginocchiaapu...](#)

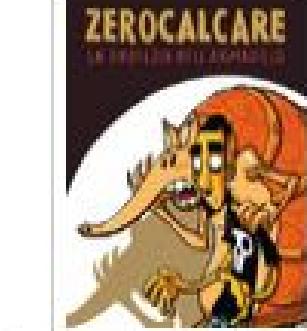
Shoutbox [See all](#)

[Leave a comment](#)

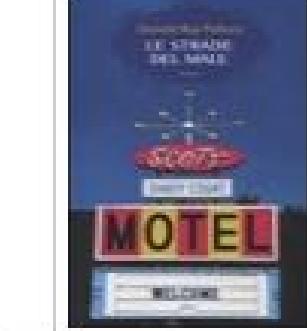
Books

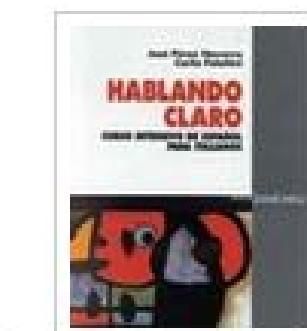
All books ▾ Rating (Highest)

 **L'Erbario delle Fate** (85)
By Benjamin Lacombe, Sébastien Perez
Finished on Nov 28, 2013 ★★★★★

 **La profezia dell'armadillo** (2931)
By Zerocalcare
Finished on Jan 6, 2013 ★★★★★

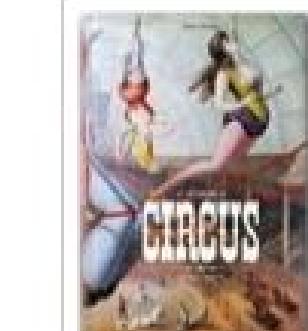
 **Biancaneve** (171)
By Jacob Grimm, Wilhelm Grimm
Finished on Feb 11, 2012 ★★★★★

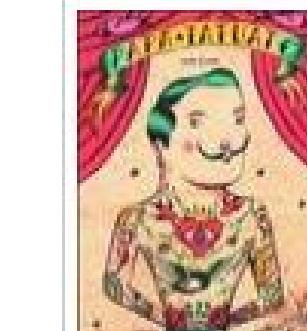
 **Le strade del male** (107)
By Donald Ray Pollock
Finished on Aug 17, 2012 ★★★★★

 **Hablando claro** (14)
Curso intensivo de Espanol para italianos
By Carla Polettini, José Pérez Navarro
Finished in 1999 ★★★★★

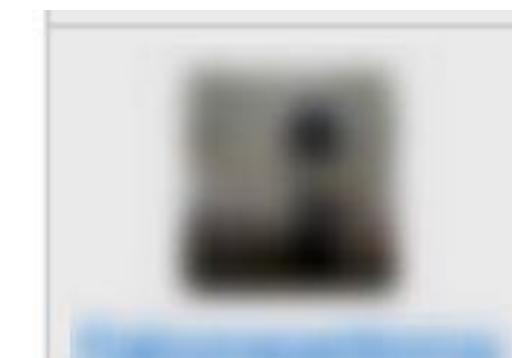
 **La ballata di Trenchmouth** (50)

 **C'era una volta...** (113)
Libro pop-up

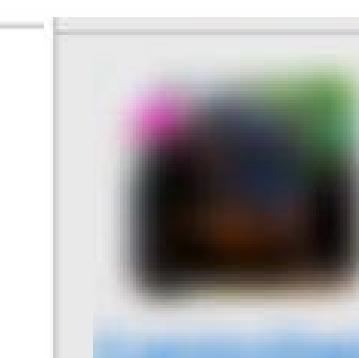
 **Circus Book 1870 - 1950** (4)
Reference

 **Papà tatuato** (67)
By Daniel

 **Venere privata** (2507)
By Giorgio



chi sei?



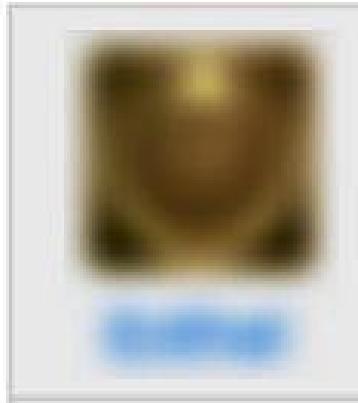
Le tue visite cominciano ad essere inquietanti....

Ieri



Grazie Lajello, mi sono divertita un sacco a leggere i commenti degli altri anobiani. Sembra un esperimento di psicologia sociale, se non ti dispiace ti aggiungo come vicino! e resisti eh...non pubblicare un libro! ;)

Due settimane fa



Ti stimo sinceramente.



LAJELLO... HAI STUFATO..NON SE NE PUO' PIU'...STA ATTENTO/A CHE SONO CAPACE DI ASSOLDARE UN HACKER PER VEDERE CHI SEI..E PO' SONO C...TUOI

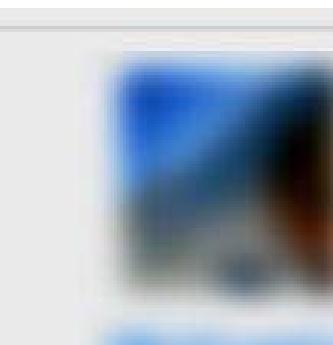
Tre settimane



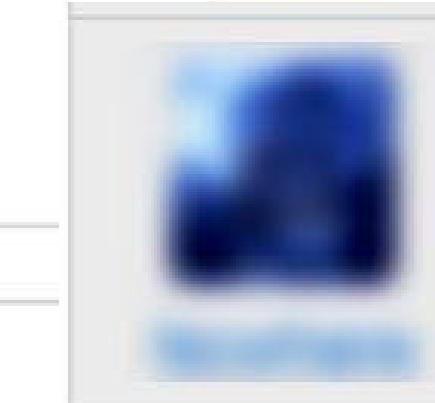
ahahahhahahaha tu sei un genio!!!!!! sei davvero un genio!!! insomma ma quante visualizzazioni hai???? sei un grande!!!! riesci a farti visitare e a farti scrivere pur non avendo libri!!! ti adoro sei grandissimo :P



ahahahaahah tu sei un genio!!

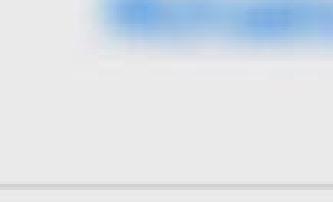


un grande.
continua così. Grazie delle visite, si vede che ti sto simpatica....



già che mi ritrovo qui mi faccio pubblicità! Venite a vedere la mia librerie è la più bella -del mondo-. (l'ultima parte andava sottolineata..)

Due settimane fa

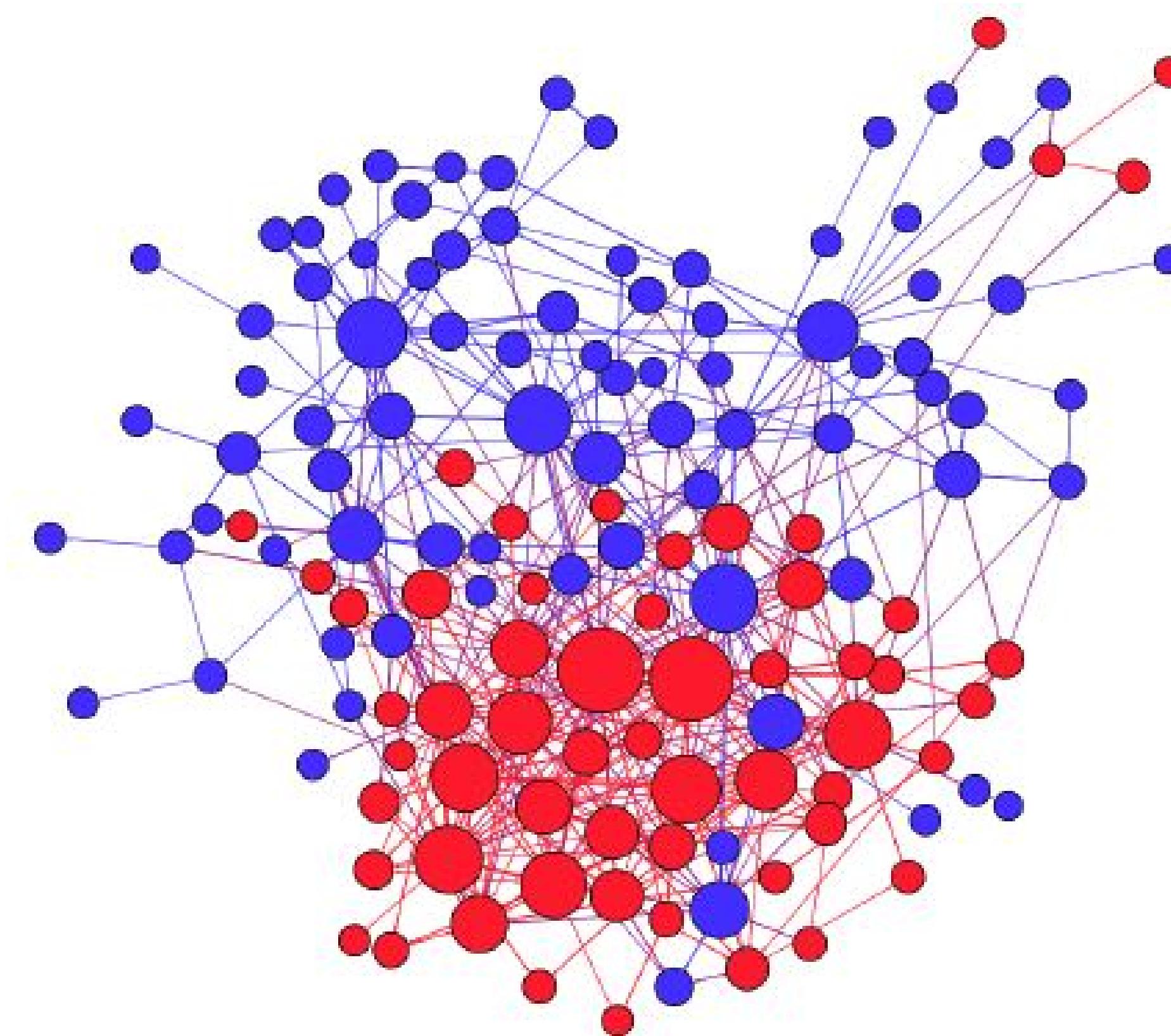


P.S: propongo di aprire un gruppo the Lajellos fans...

3 giorni fa



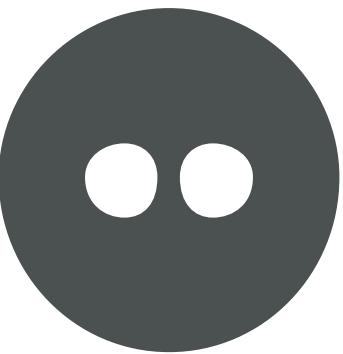
From Influence to Chaos



2.

QUALITY INSPIRES QUALITY

DATA



15B

PHOTOS

40M

USERS

550+M

SOCIAL TIES

10+

YEARS

GROUPS

FAVORITES

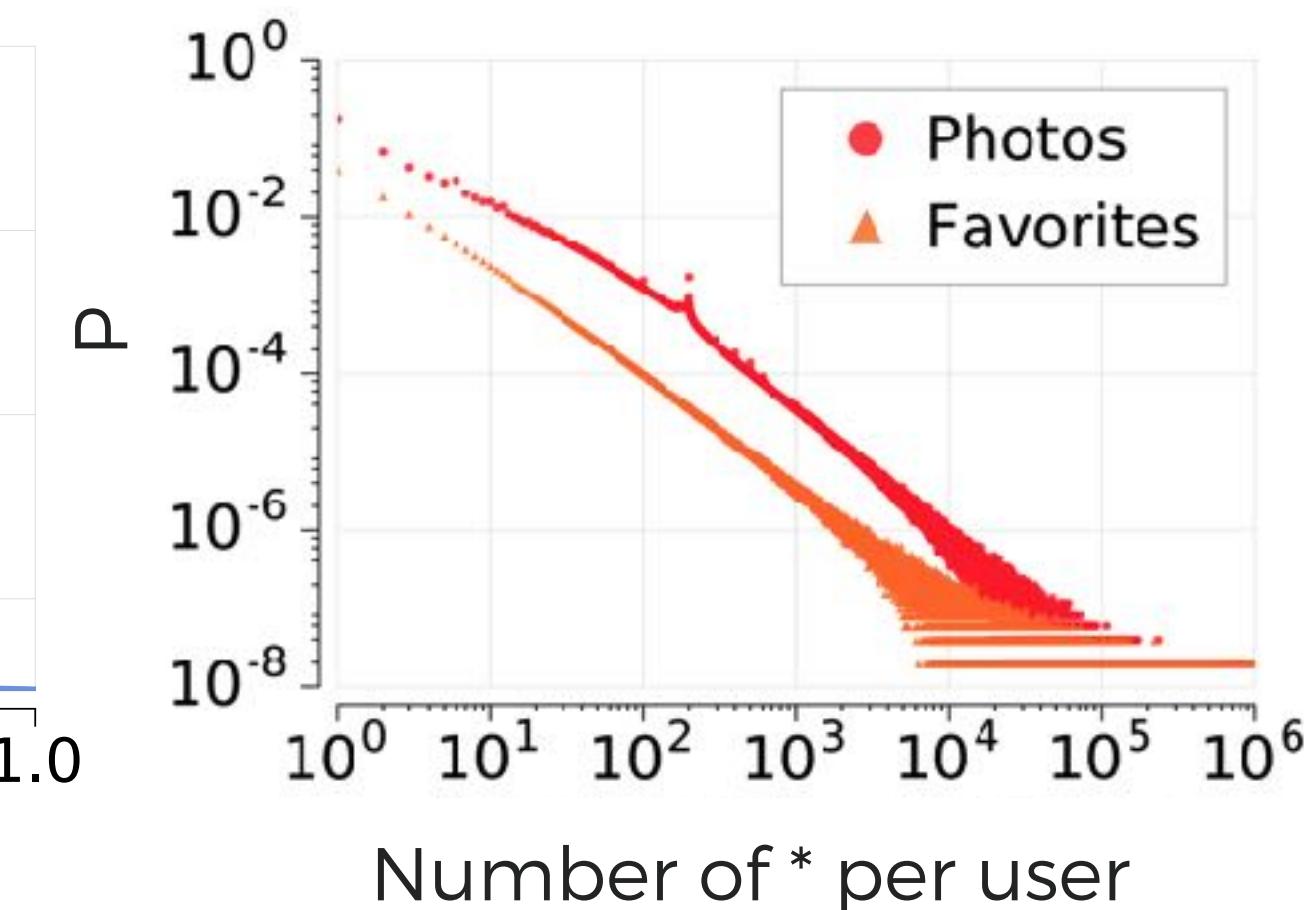
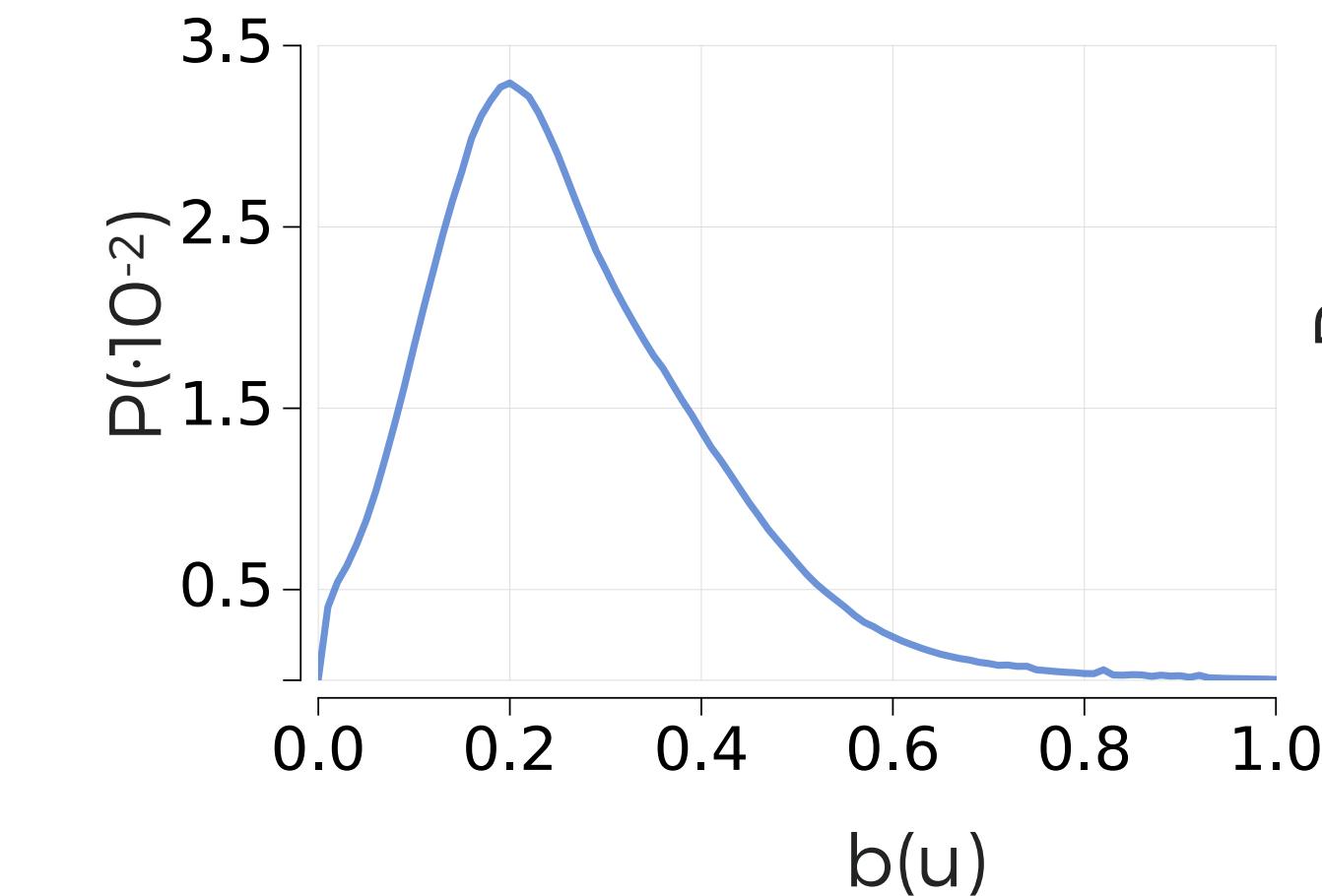
Quality at user-level

$b(u)$: average of beauty scores

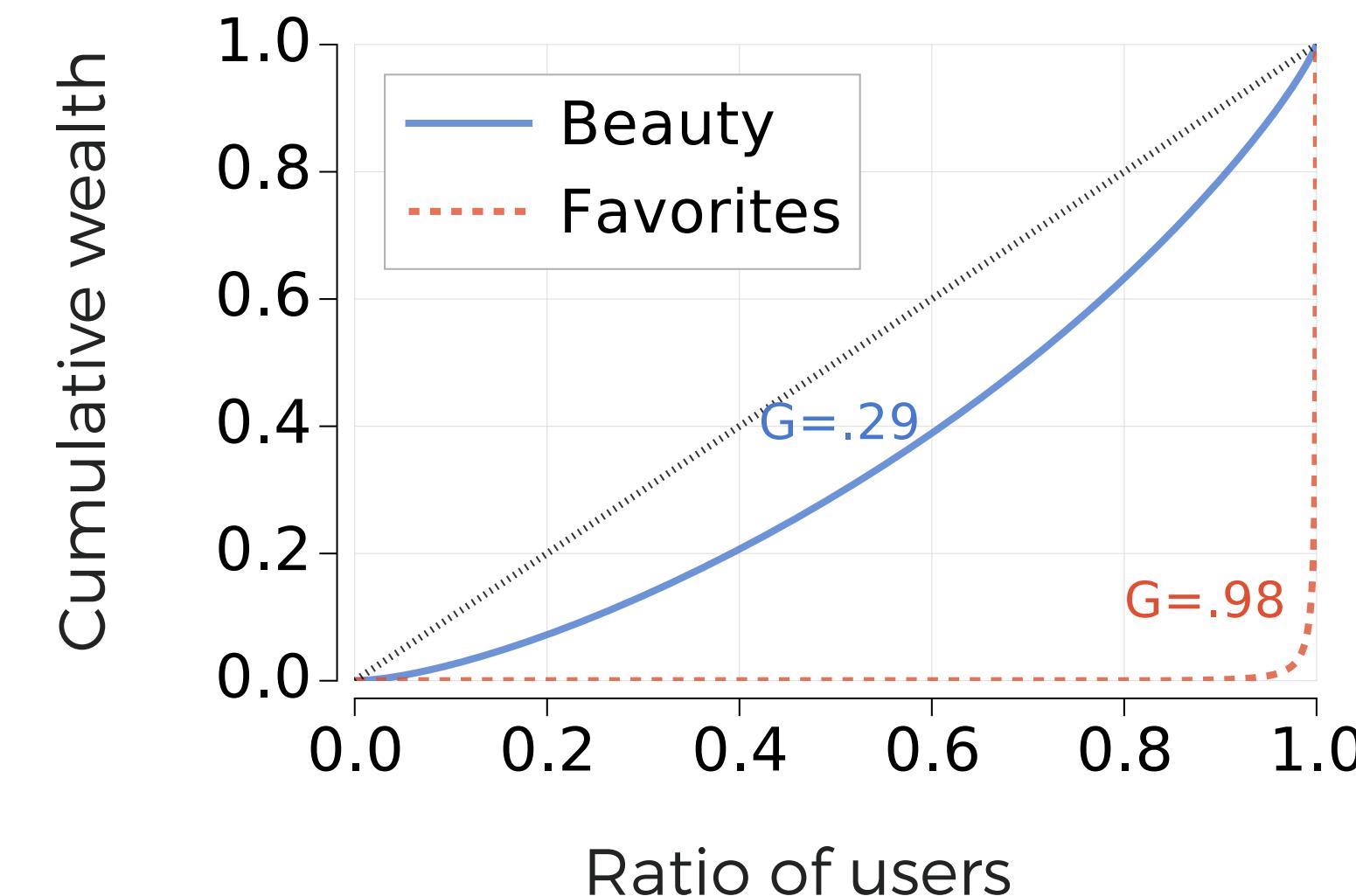
we are currently working on
different strategies
evolution of photographers
careers

**Quality distribution is highly
unequal ($G=0.98$)**

measuring beauty at scale (15B pics)



inequality (Gini's coefficient)



QUALITY AND CONNECTIVITY ARE (WEAKLY) CORRELATED

quality(u) is

in-degree

$(\rho = .22)$

out-degree

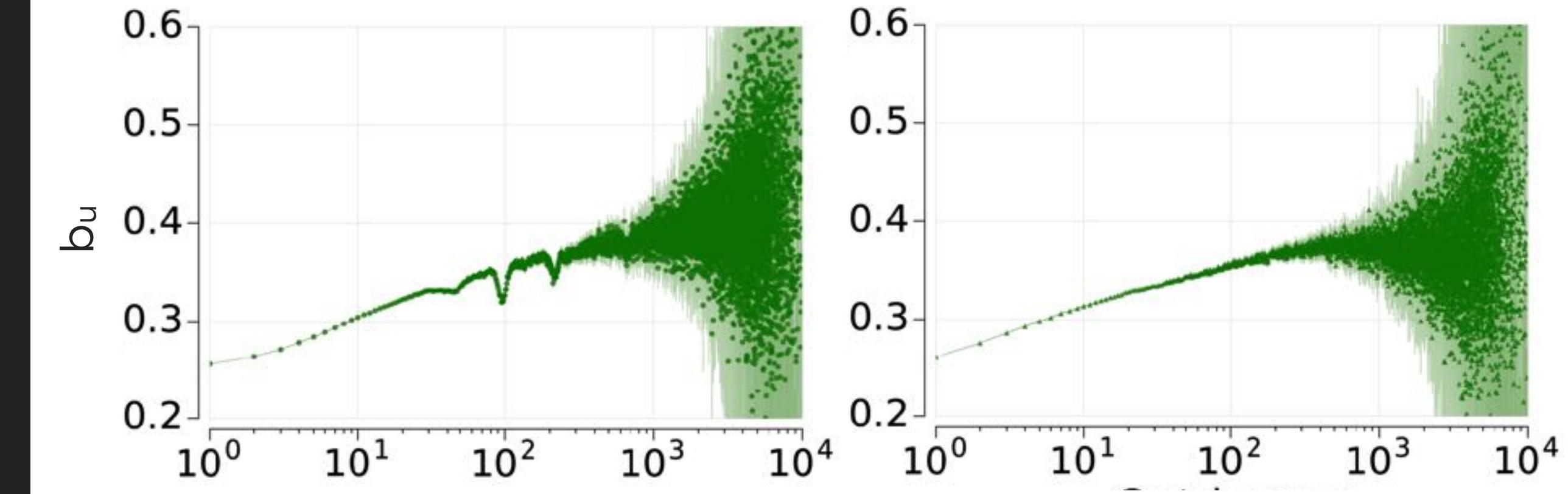
$(\rho = .24)$

#favorites

$(\rho = .17)$

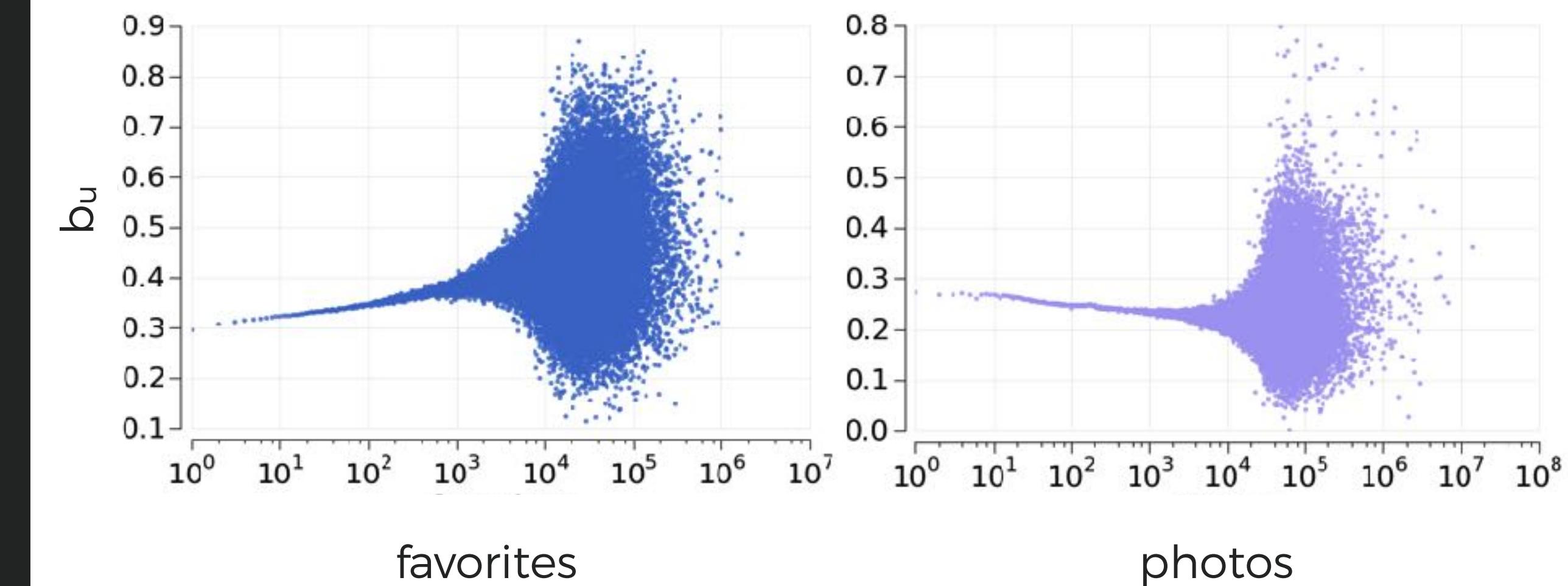
#photos

$(\rho = -.03)$



in-degree

out-degree



favorites

photos

SOCIAL TIES LINK USERS WITH SIMILAR QUALITY

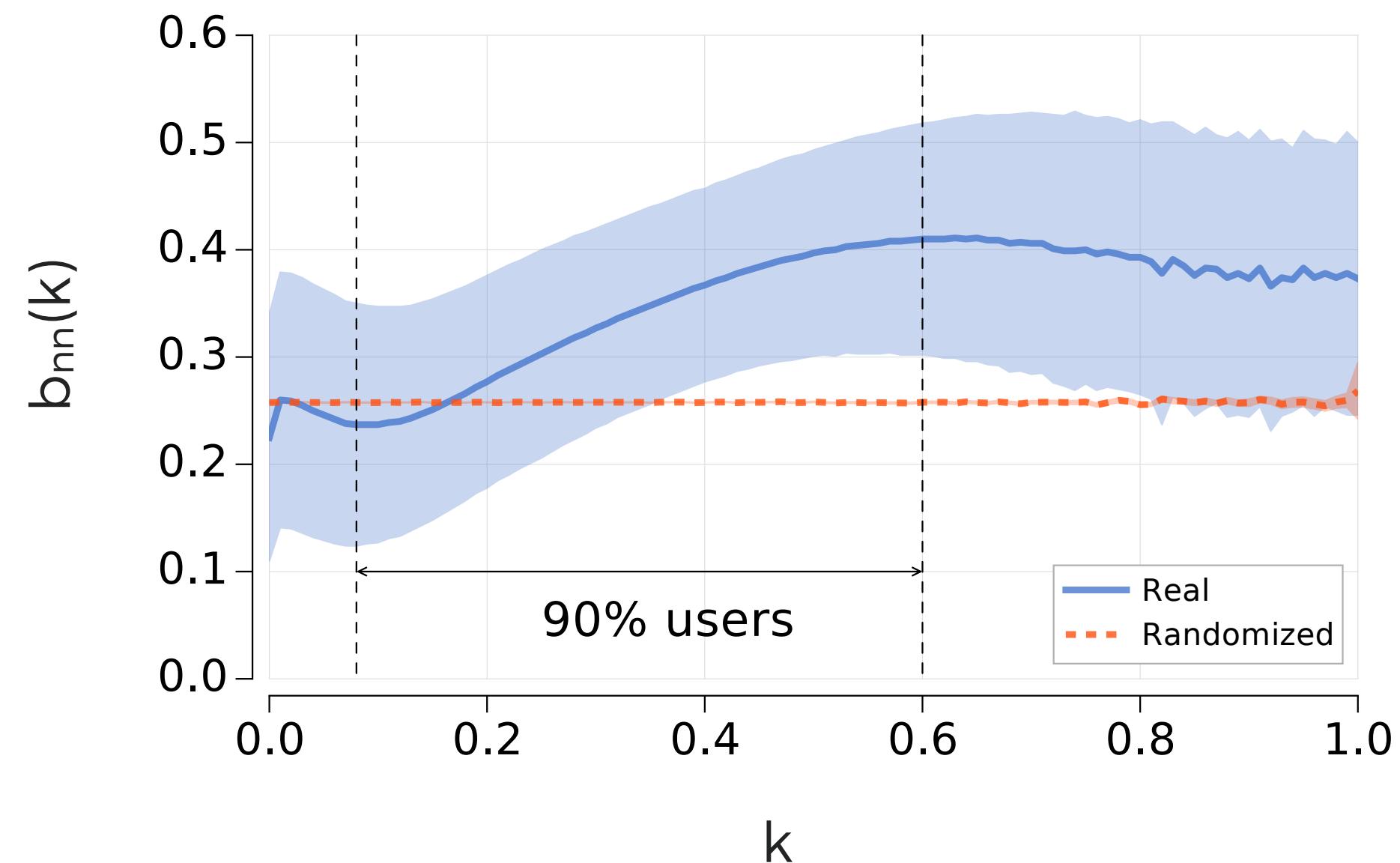
assortative mixing

users tend to be linked to accounts that publish photos with similar quality

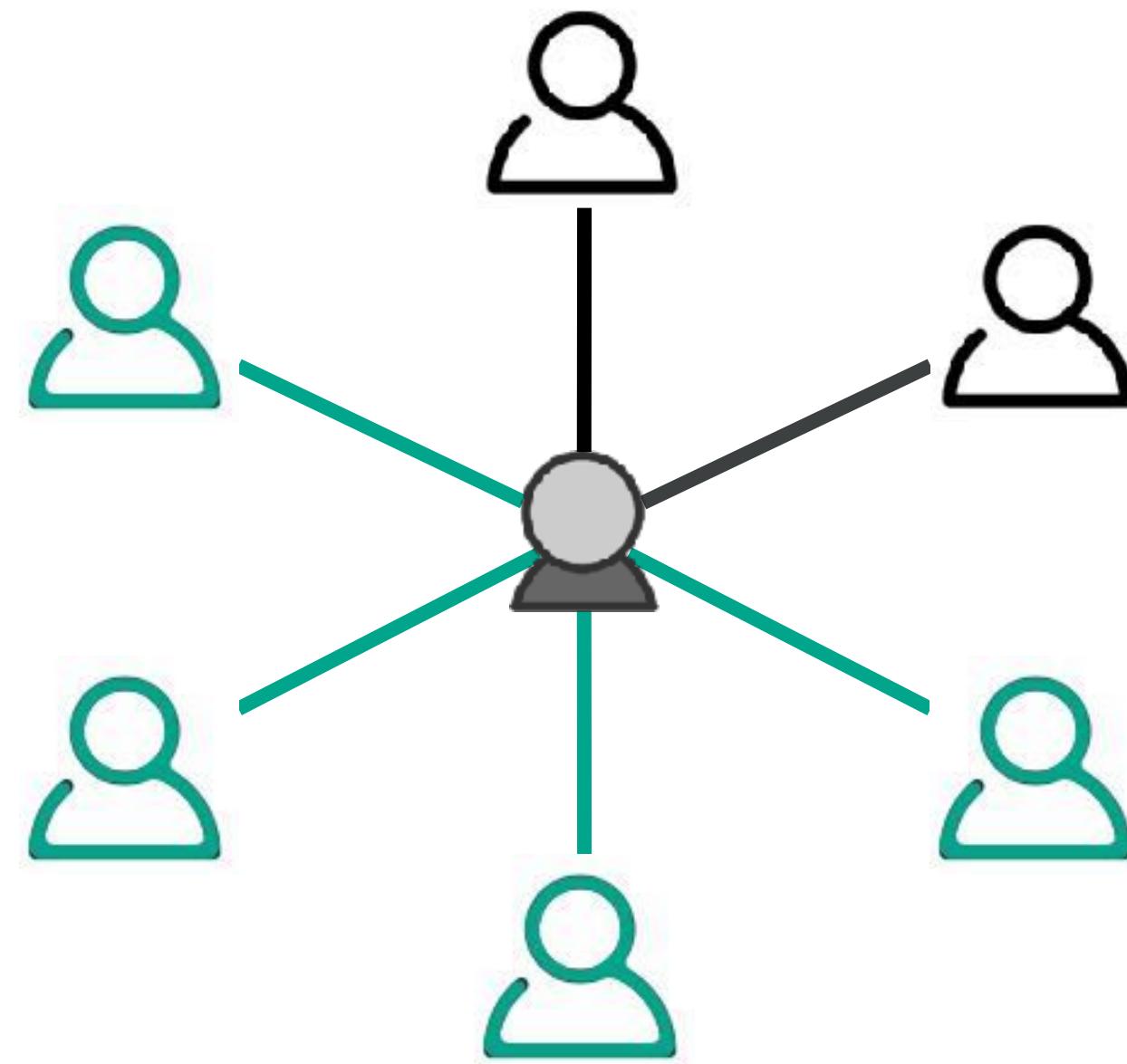
(Spearman $\rho = 0.48$)

This could be ascribed to homophily influence

$$b_{nn}(k) = \frac{1}{|\{i : \bar{b}(i) = k\}|} \cdot \sum_{j \in \Gamma_{out}(i), \forall i; \bar{b}(i)=k} \bar{b}(j)$$



Majority Illusion



0.26
 μ

43%
above μ

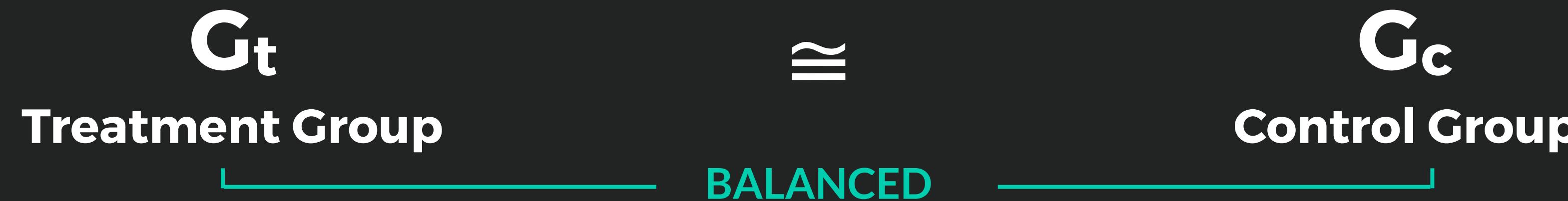
65%
 $> 43\% \text{ friends above } \mu$

To summarise

- High unequal distribution of quality
- Quality and connectivity are weakly correlated
 - no correlation with activity
- Users tend to connect to similar quality users
- Local perception of higher quality of neighbours

Matching experiments

Inferring causality from observational data [Rubin 2001][Stuart 2010]
[Althoff et al. 2017]



Two groups are **balanced** on a **covariate X** when the **standardized bias SB_X** is under a given threshold (usually set to 0.25)

$$SB_X(G_t, G_c) = \frac{\bar{X}_t - \bar{X}_c}{\sigma(X_t)} \leq 0.25$$

USER

in-out-degree

photos

groups

favorites

 μ (neighs) μ (beauty)

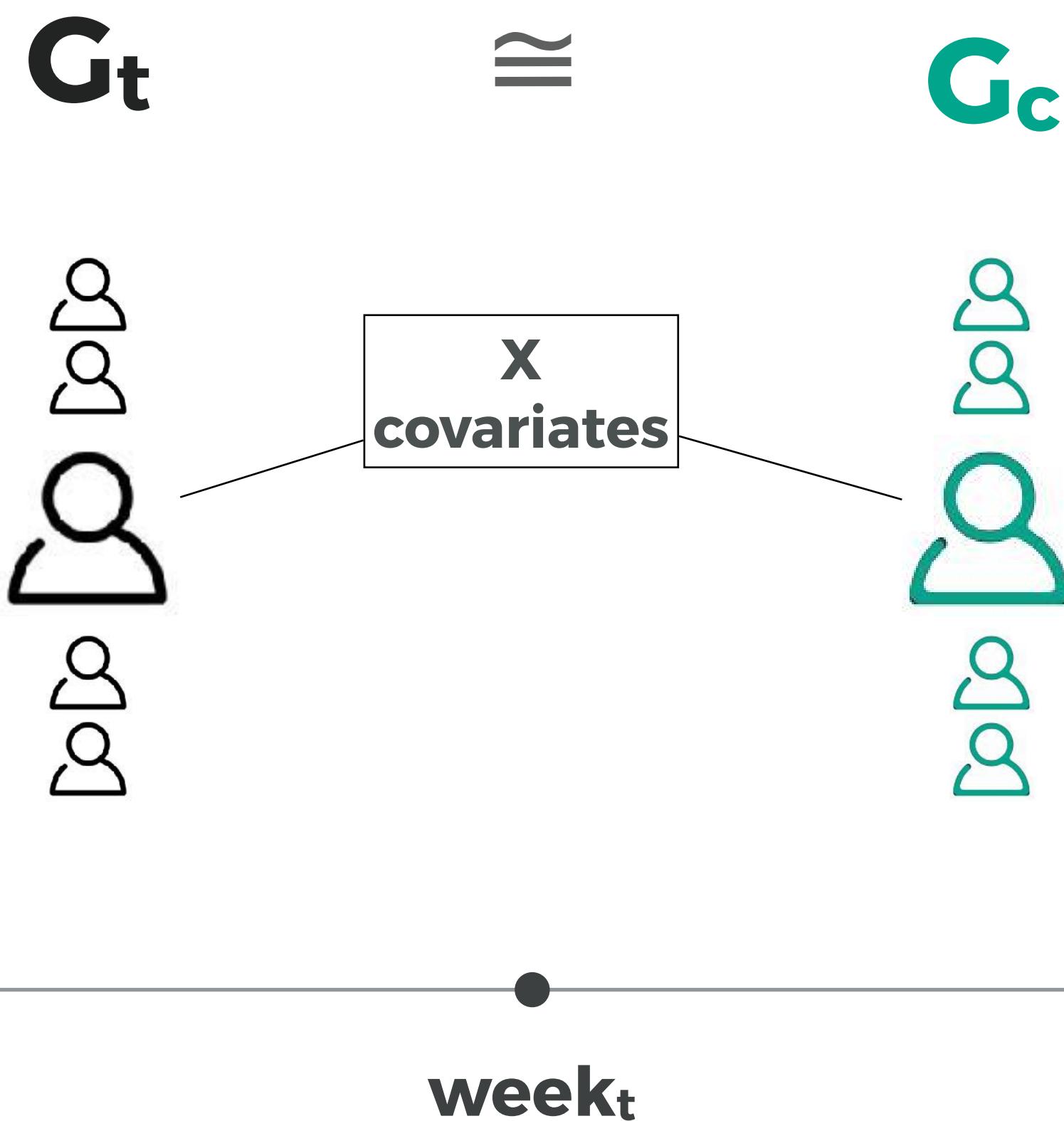
weeks from join date

NEIGHBORS

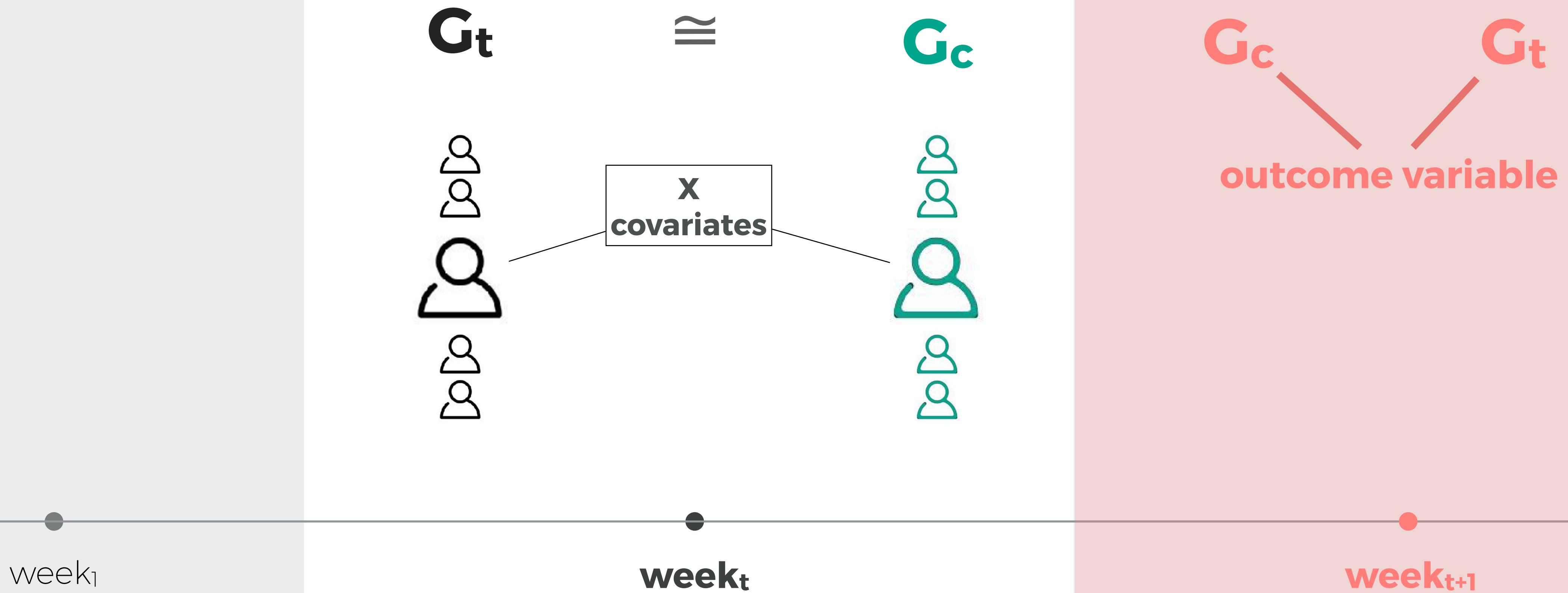
photos

 μ (beauty)

Setup on Flickr data



Setup on Flickr data



QUALITY INSPIRES QUALITY

treatment:

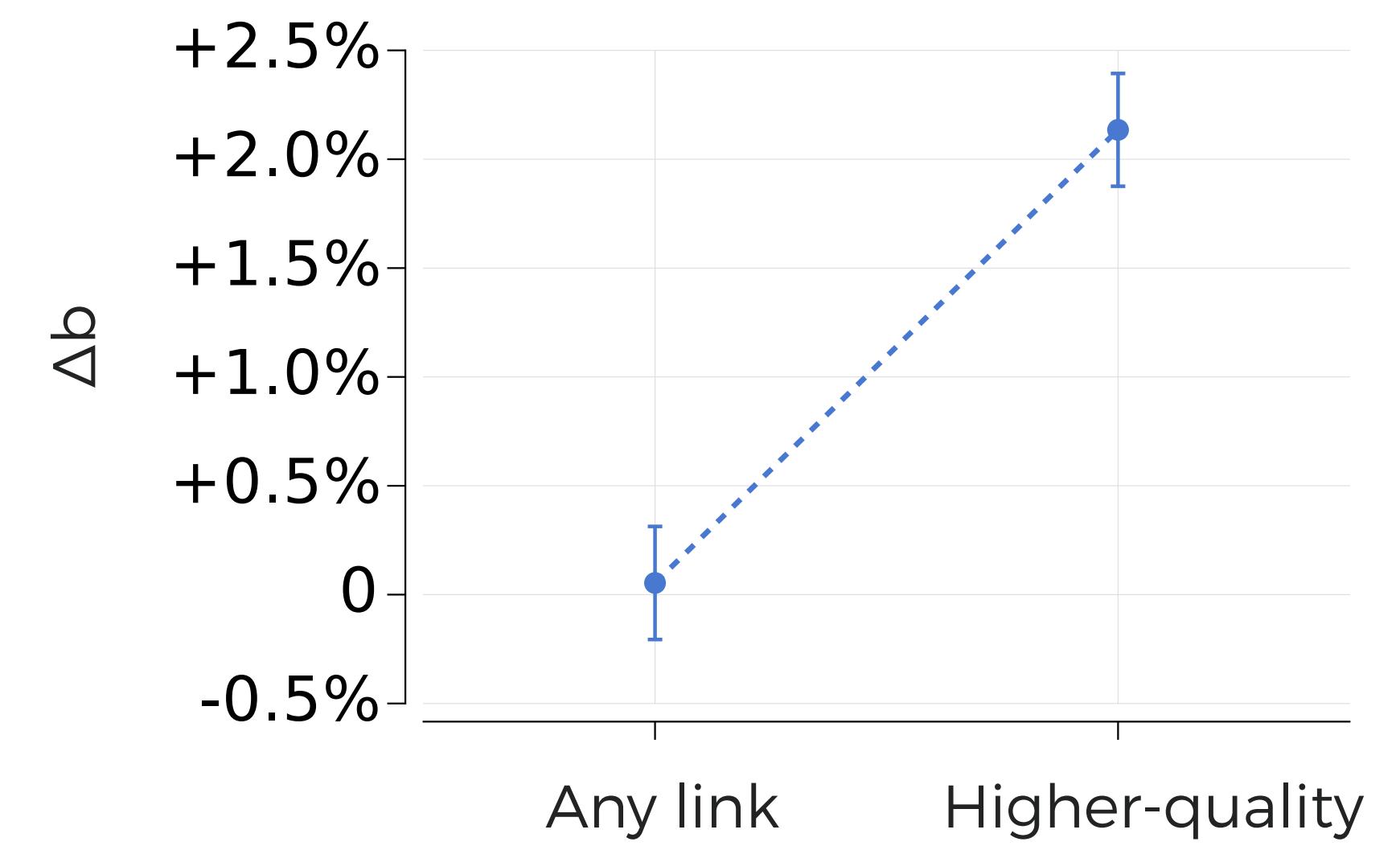
new tie to a **higher-quality** neighbor

outcome variable:

increase in quality in future photos (i.e., photos made in w_{t+1})

EIGHBOURS QUALITY AFFECT USER QUALITY

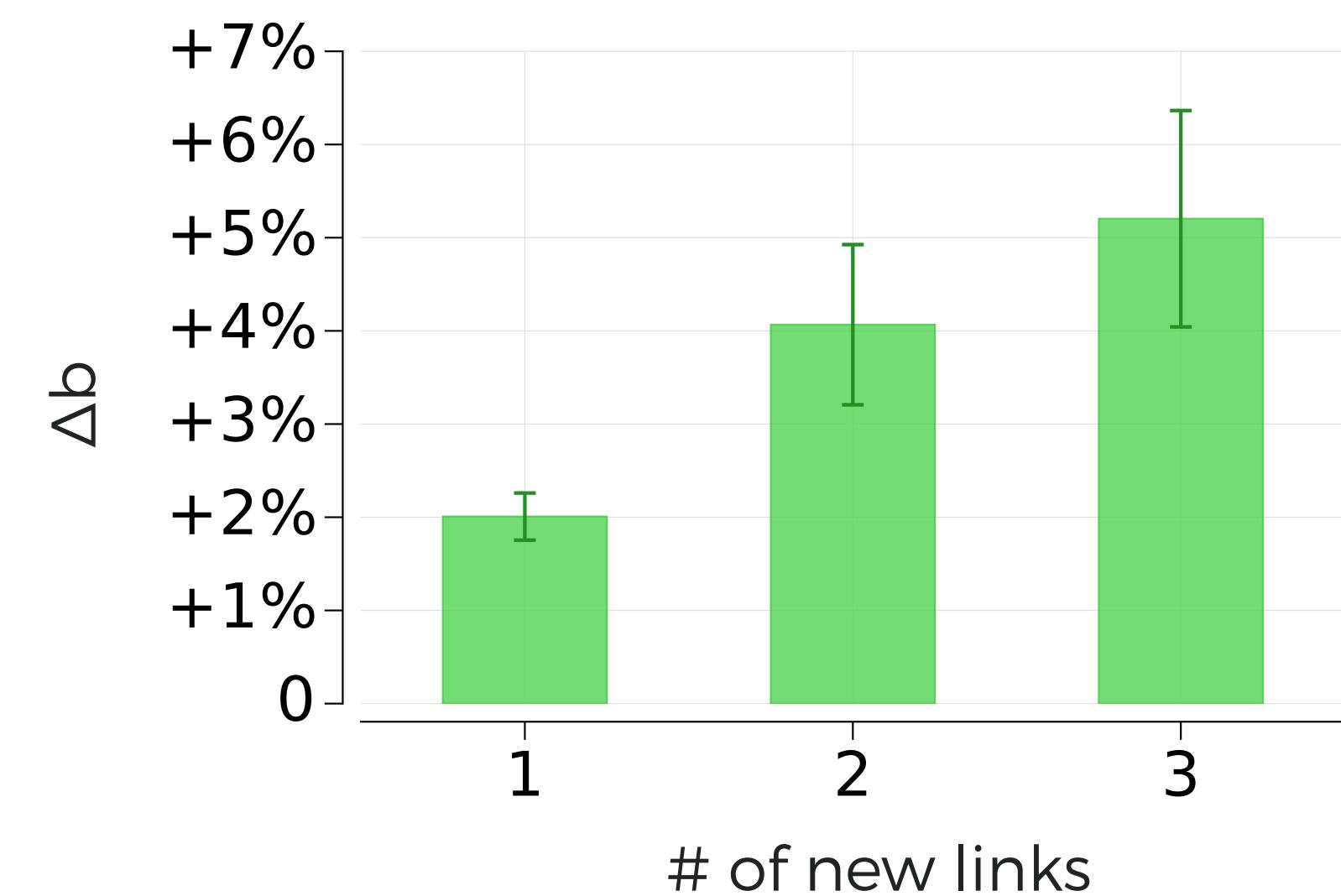
1) **G_t experiences an average 2% increase, no significant increase for G_c**



EIGHBOURS QUALITY AFFECT USER QUALITY

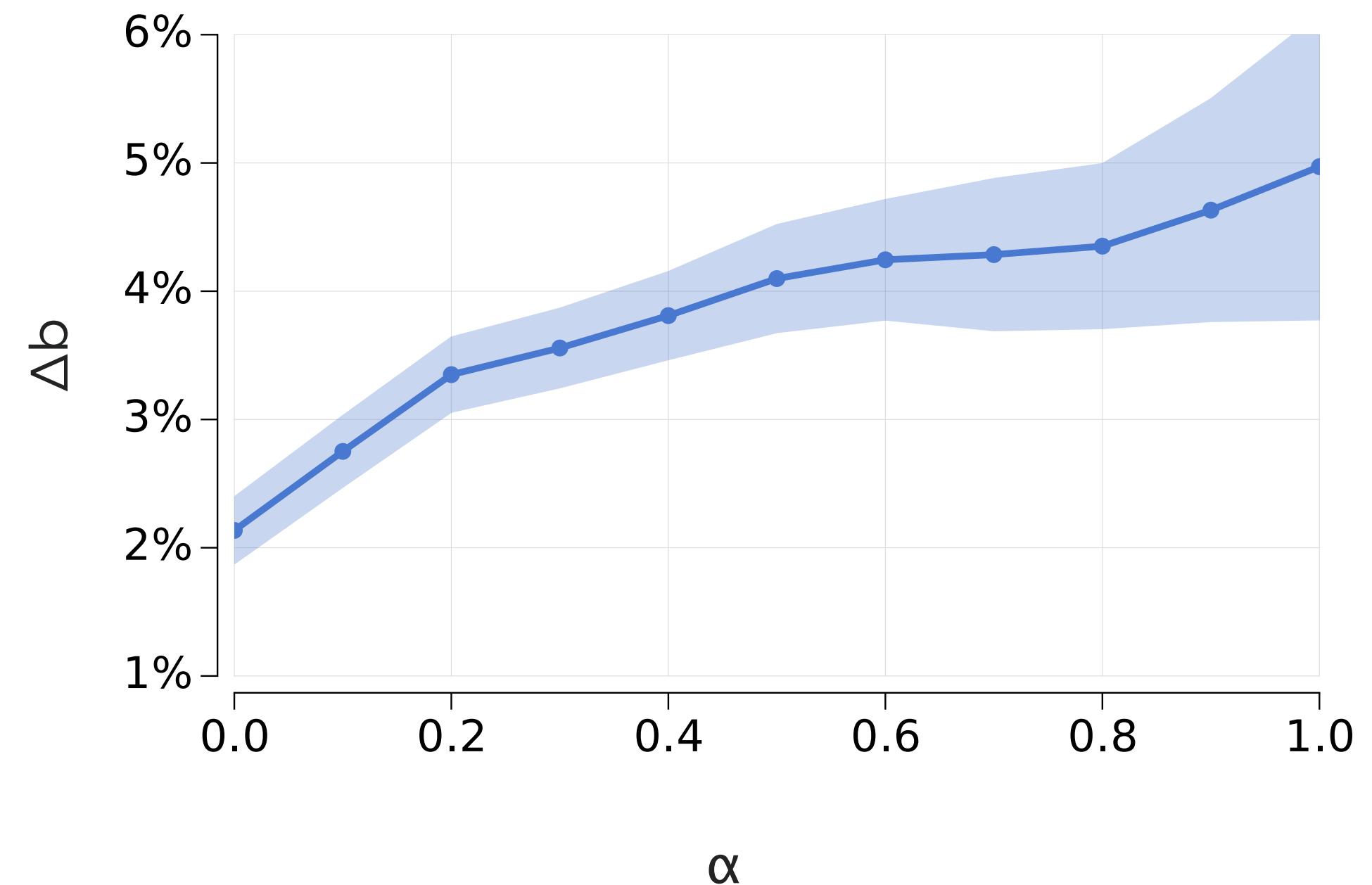
1) G_t experiences an average 2% increase, no significant increase for G_c

2) Influence effect accumulates with new connections, with diminishing returns



EIGHBOURS QUALITY AFFECT USER QUALITY

- 1) G_t experiences an average 2% increase, no significant increase for G_c**
- 2) Influence effect accumulates with new connections, with diminishing returns**
- 3) The greater the beauty differential, the greater the increase, noticeable until $\alpha=0.5$**



IMBALANCE DROPS ENGAGEMENT

treatment:

average neighbor beauty deviates from the user beauty (**high imbalance**)

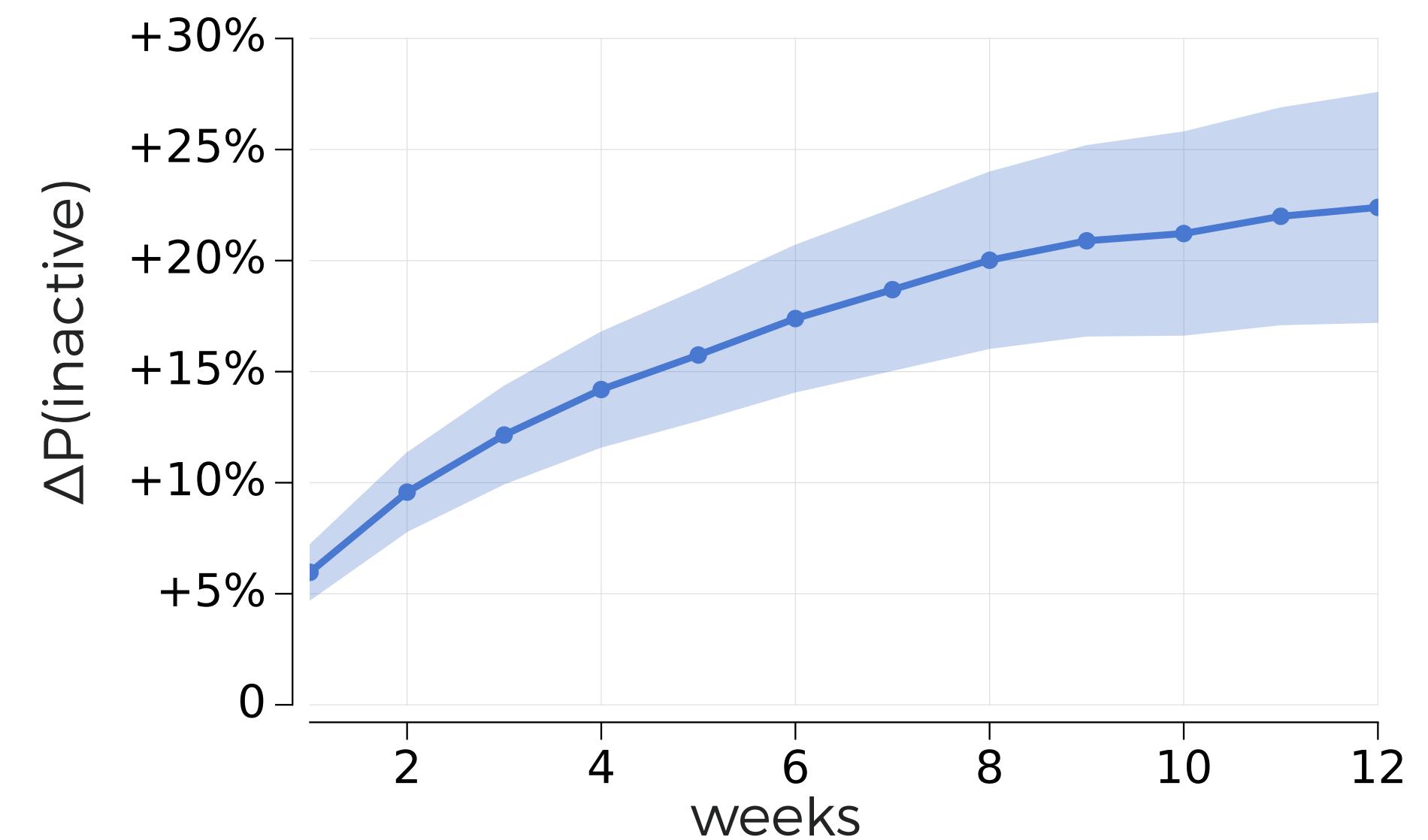
outcome variable:

proportion of users in each group who remain **inactive** (i.e., no photo uploads)

QUALITY IMBALANCE AFFECT ENGAGEMENT

G_t has higher probability of inactivity than grows from +5% to +20% in the first 12 weeks

People exposed to photos that deviate significantly, in terms of quality, from their own contributed content, are more likely to become disengaged in the future



The Social Brain Theory

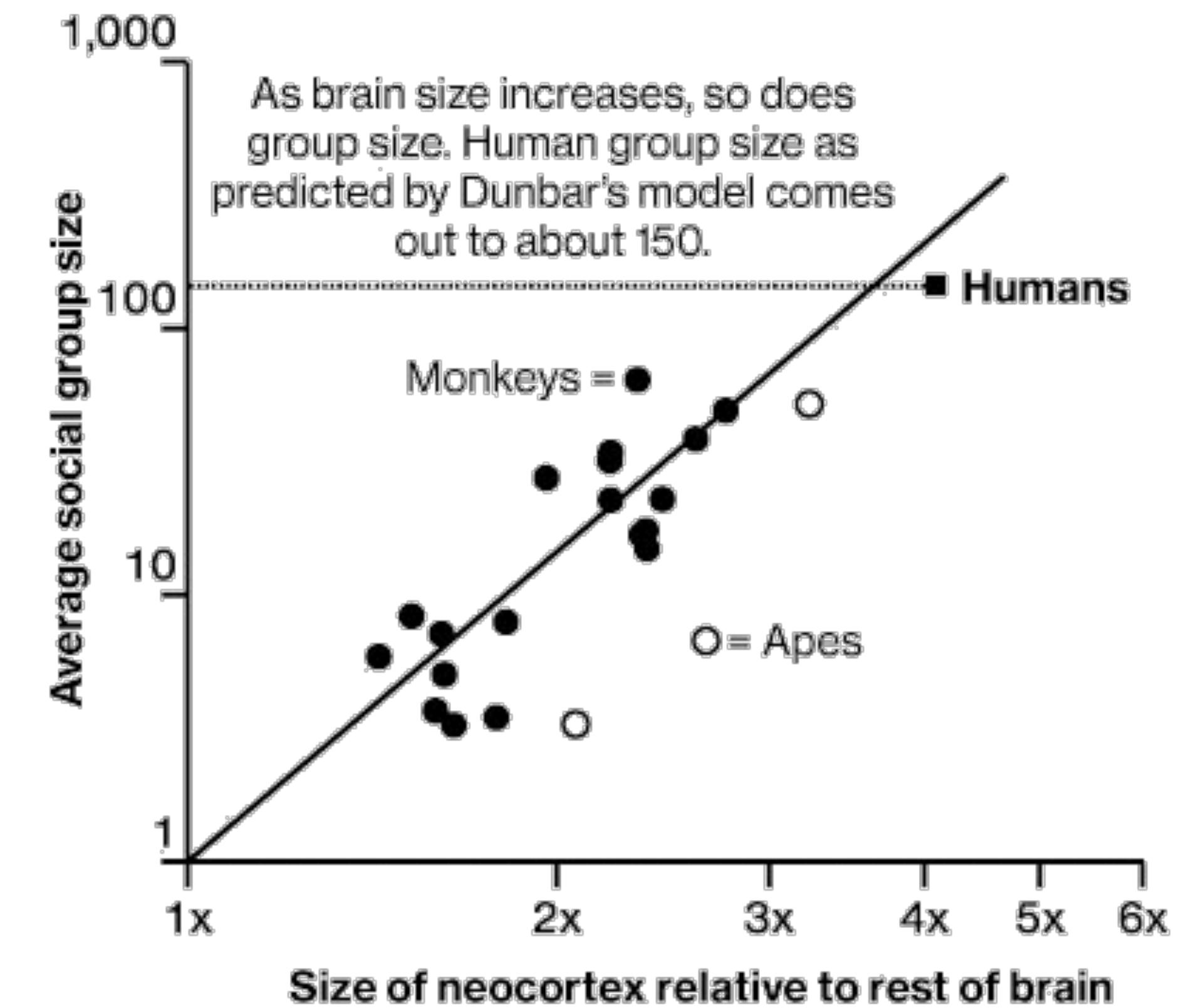
Densification of the Social Graph

- **Why does the triangle closure process stops?**
 - Time and space constraints
 - Attention and cognitive boundaries
- **How big is a person's ego network?**
- **Is there a universal limit to human connectivity?**

The Social Brain Theory

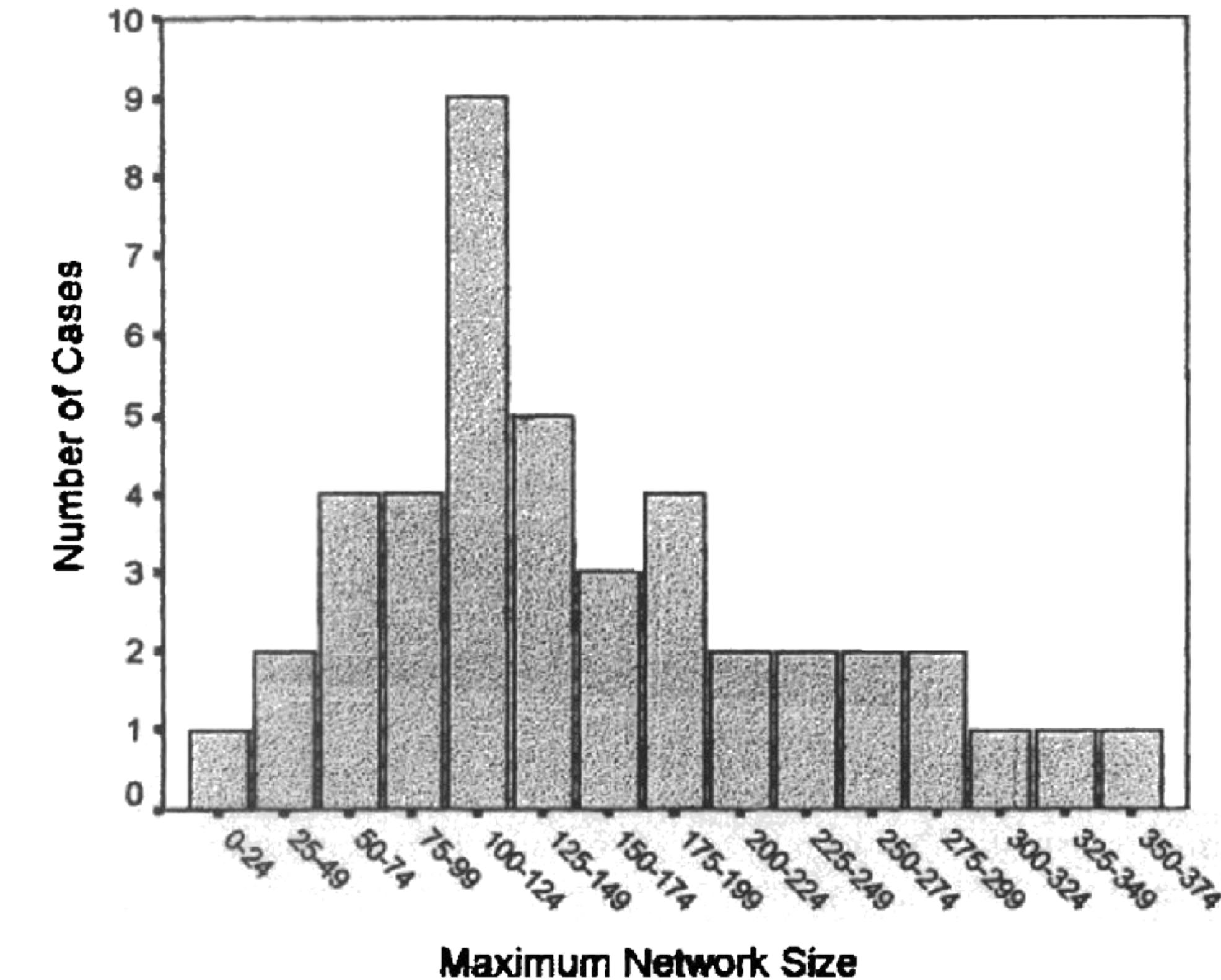
- **Big brains evolved to solve the problem of social life**
 - Cognitive effort
 - Time constraints (social grooming)
 - Face to face interaction
 - Emotional intensity
- **Correlation between neocortical volume and typical social group size**
- **Limits of human social network size between 100 and 200 individuals** (150 is the anecdotal number – Dunbar number)

The Social Cortex



Exchange of Xmas cards

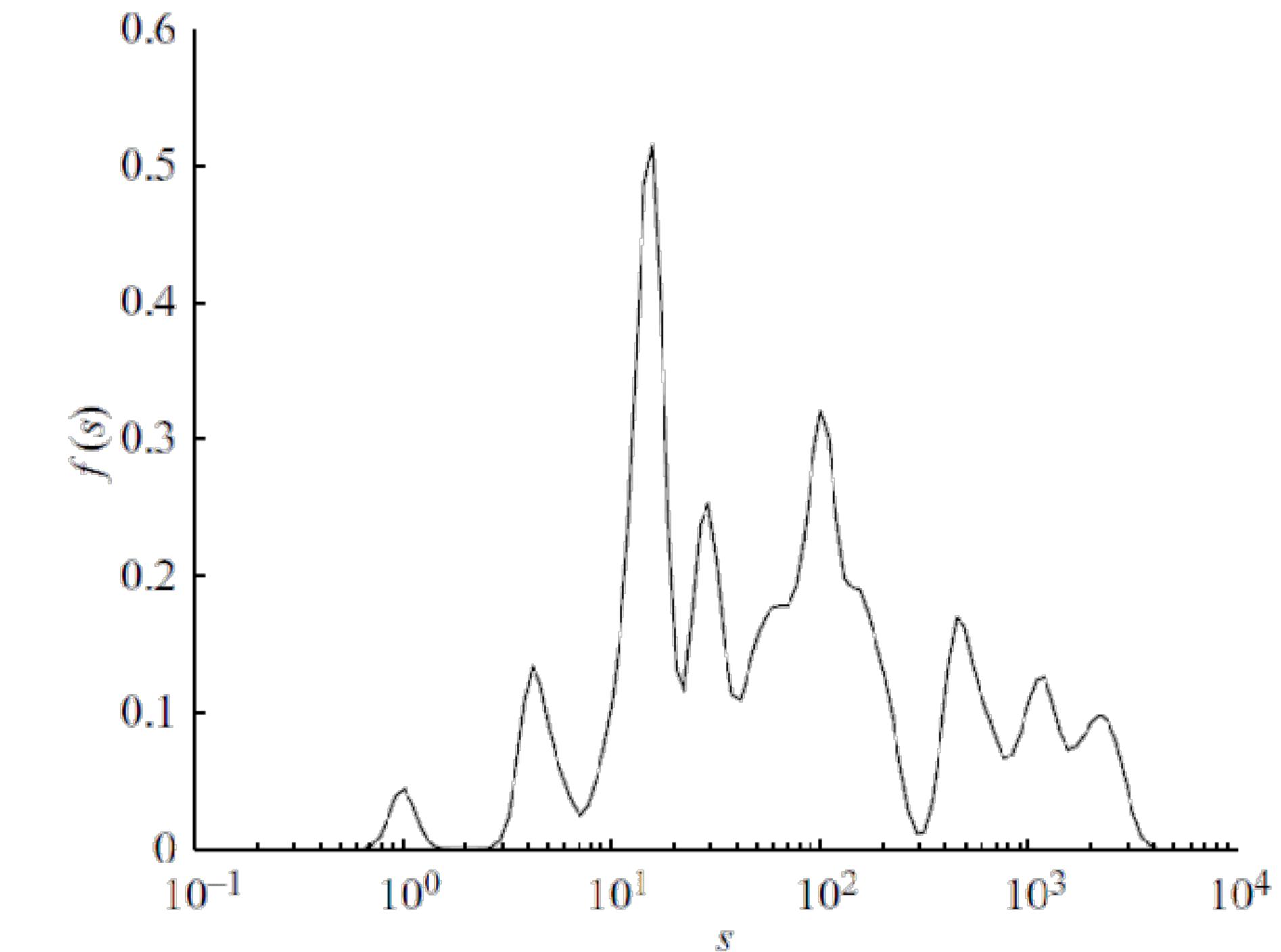
- **Sending Christmas cards individuals make an effort to contact their network of acquaintances they value**
- 43 questionnaires from UK households
- 2,984 Christmas cards sent
- Number of contacts =
 - # family members +
 - # card recipients (household) +
 - # people seen on Christmas
- Results
 - Mean cards = 68.19
 - **Mean contacts = 153.5**



Experimental limitations ?

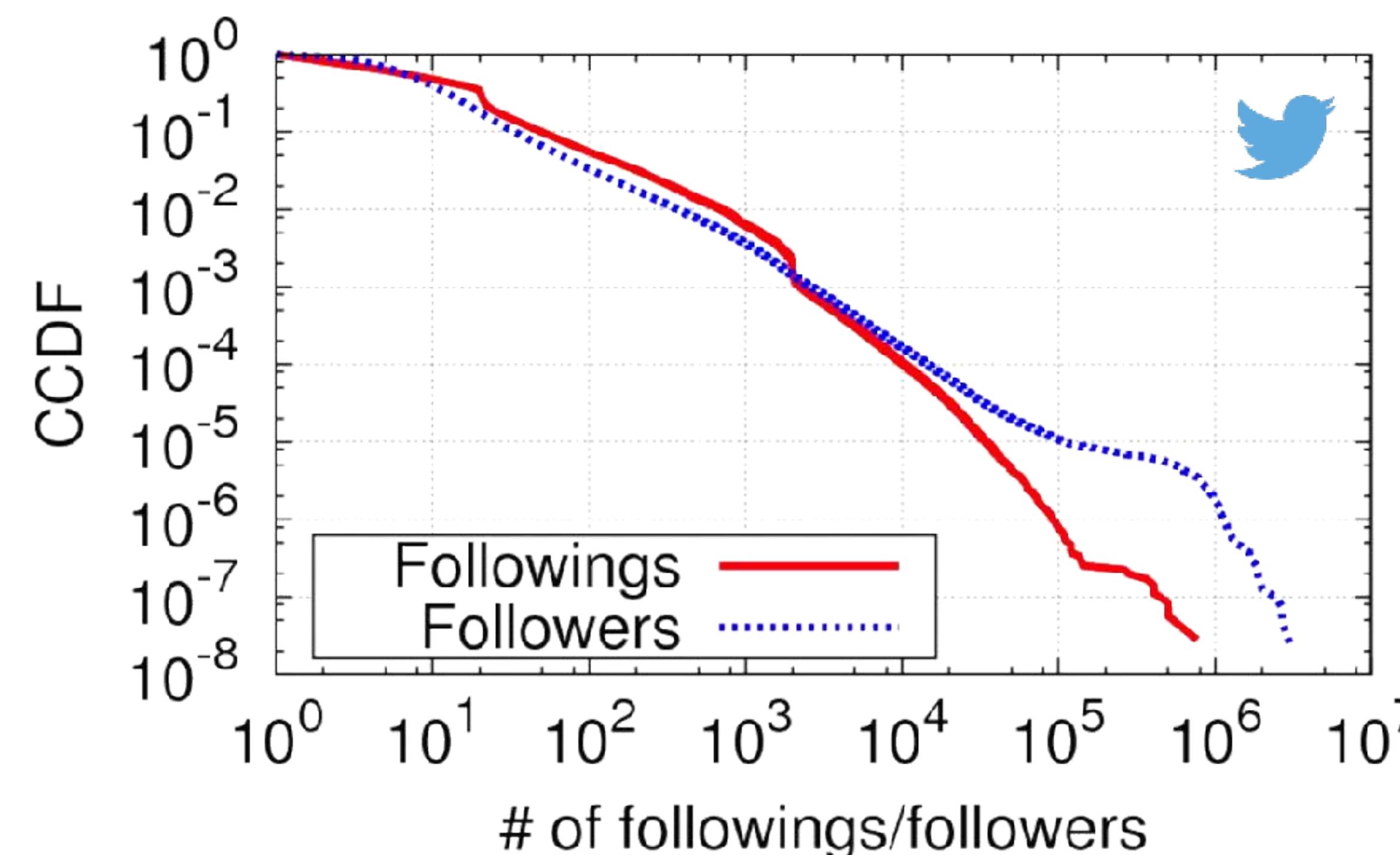
Social Circles

- **61 group datasets from literature**
- **Expanding intimacy circles**
 - Support cliques (5)
 - Sympathy circle (15)
 - Close relationships (50)
 - Stable relationships (150)
 - Acquaintances (500)
 - People we can name (1500)
- **Scaling factor of ~3**



Grooming in digital media?

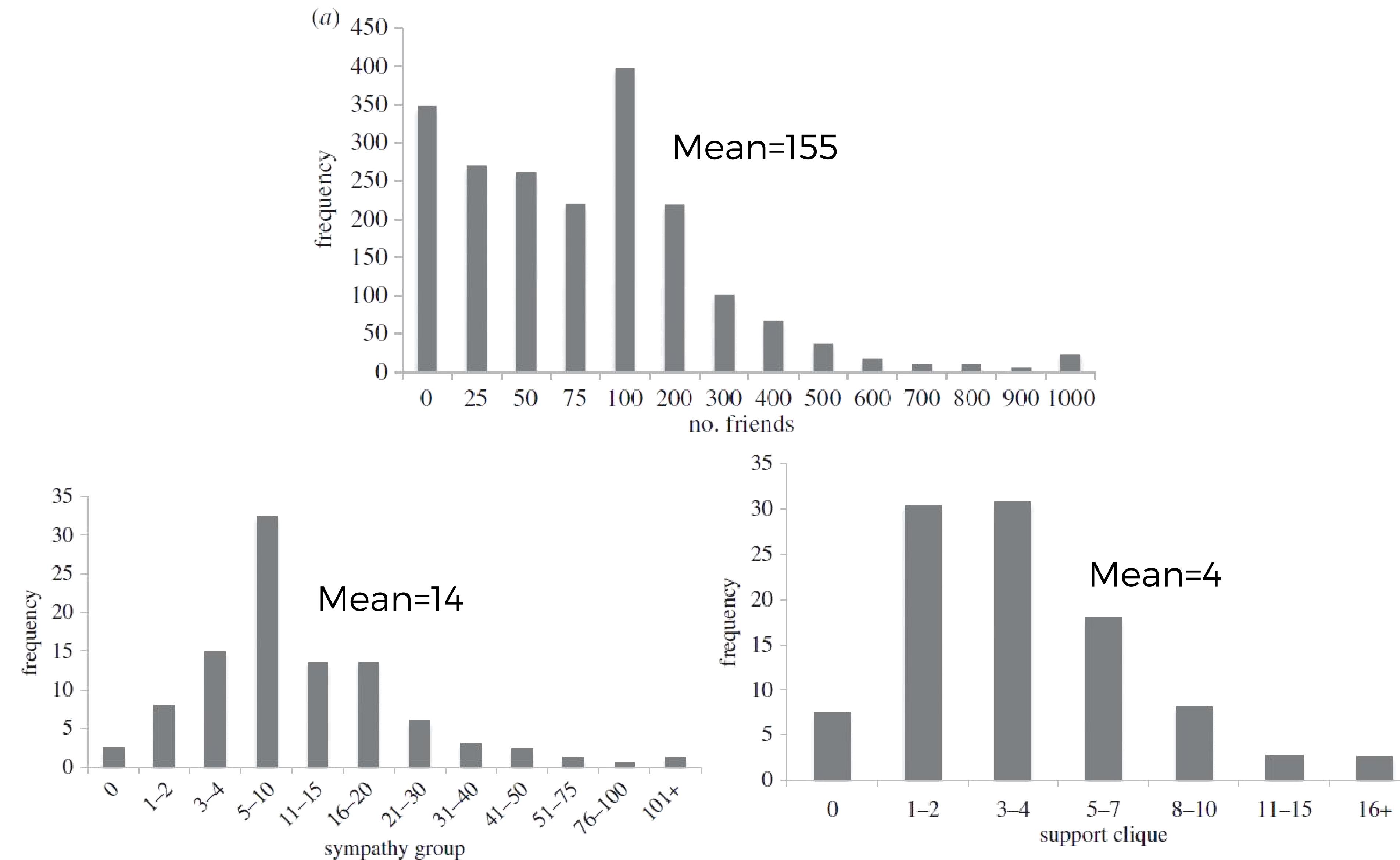
- SNSs may circumvented at least some of the real-life constraints
- Apparently no limit ... at least for some people



Probing digital barriers

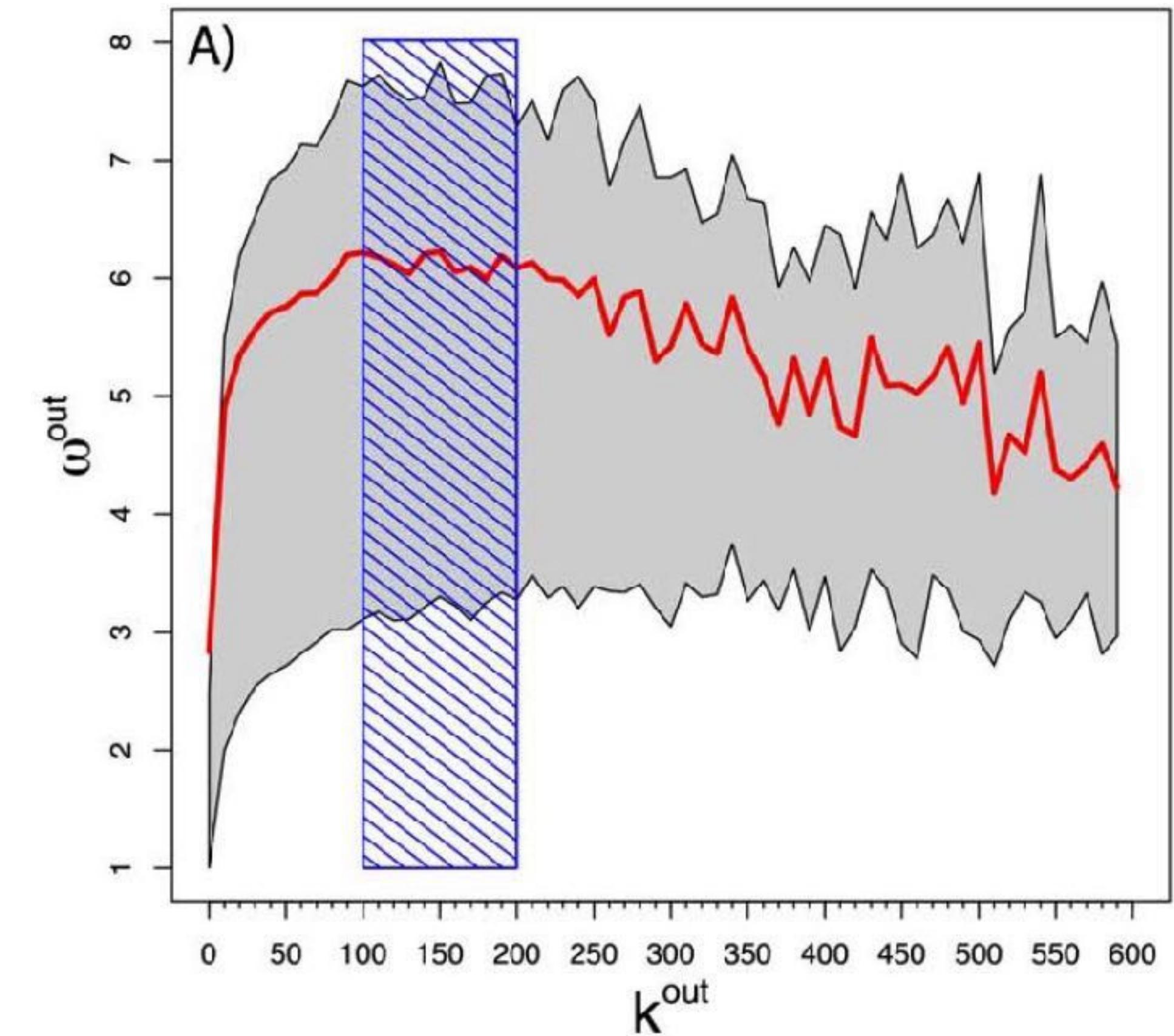
- **Survey to 2,000 UK Facebook users**, stratified sample (age, gender and other cofounding)
- How many Facebook...
 - **Friends [0, 1,000+]**
 - **Close friends [0, 100+]**
 - Support clique members (advice/sympathy in time of distress) [0, 16+]

SNSs do not circumvent real-life constraints



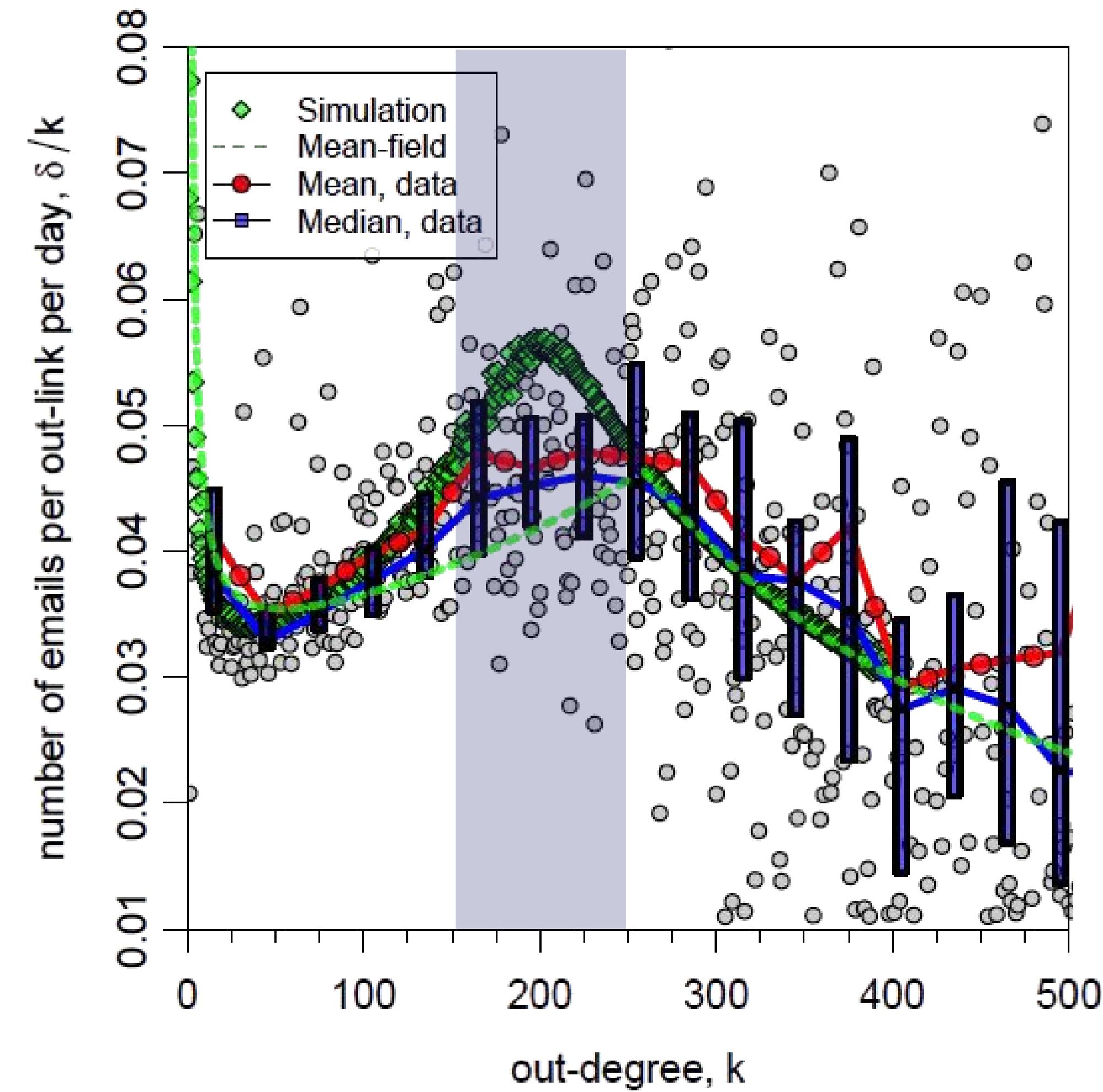
Dunbar number on Twitter

- **Twitter dataset, 380M tweets, 1.7M users**
- **Network of conversations**
- Direct measurement of the interaction strength:
$$\omega_i^{\text{out}}(T) \equiv \frac{\sum_j w_{ij}(T)}{k_i^{\text{out}}}$$
- w_{ij} number of replies
- T time window
- K_{out} reciprocated connections
- Plot ω^{out} (av. link strength) vs. k^{out} (number of bi-links)



Dunbar number on email

- Email corpus of Univ. Oslo employees and students (35,600)
- ~1M contacts in the outside world



Group perspective

- Flickr dataset, 238M comments, 504K groups, 71M users
- Measure activity inside the group:

$$a_g = \frac{E_g^{\text{int}}}{(D_g^{\text{in}} D_g^{\text{out}})/E}$$

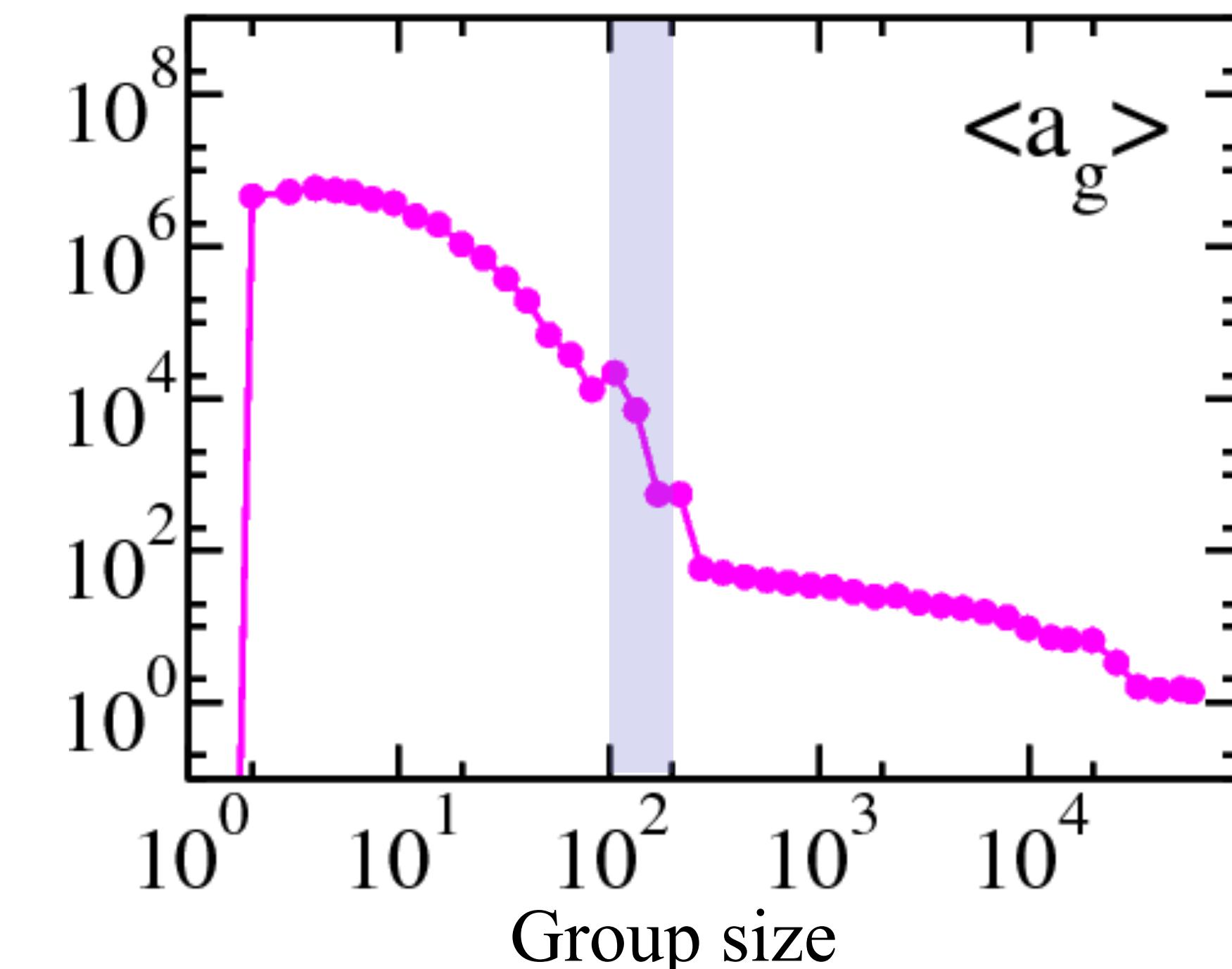
E=tot #edges

E_g^{int} =#internal group edges

D_g^{in} =#edges incoming to group

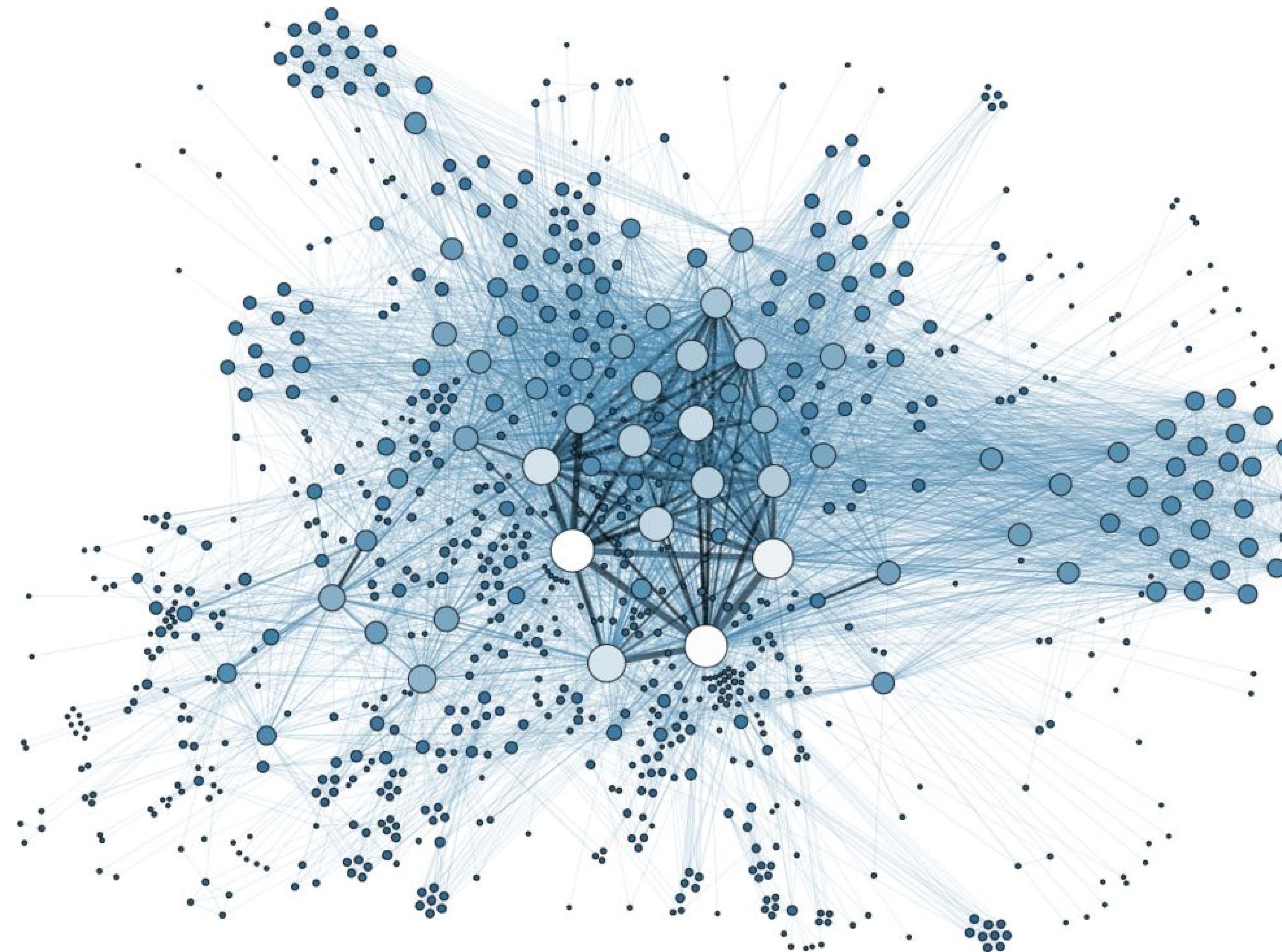
D_g^{out} =#edges outgoing from group

- Plot avg. a_g vs. group size

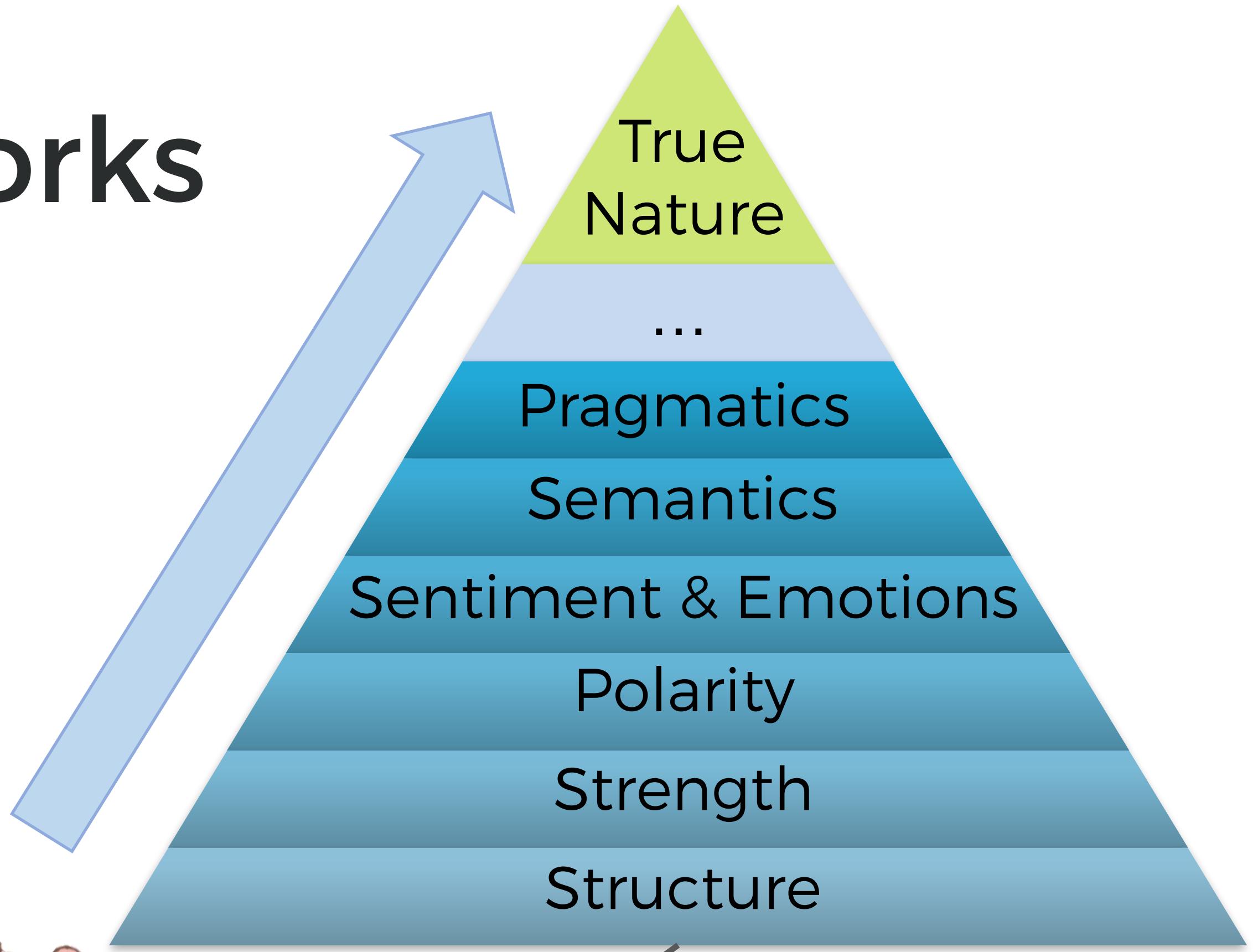
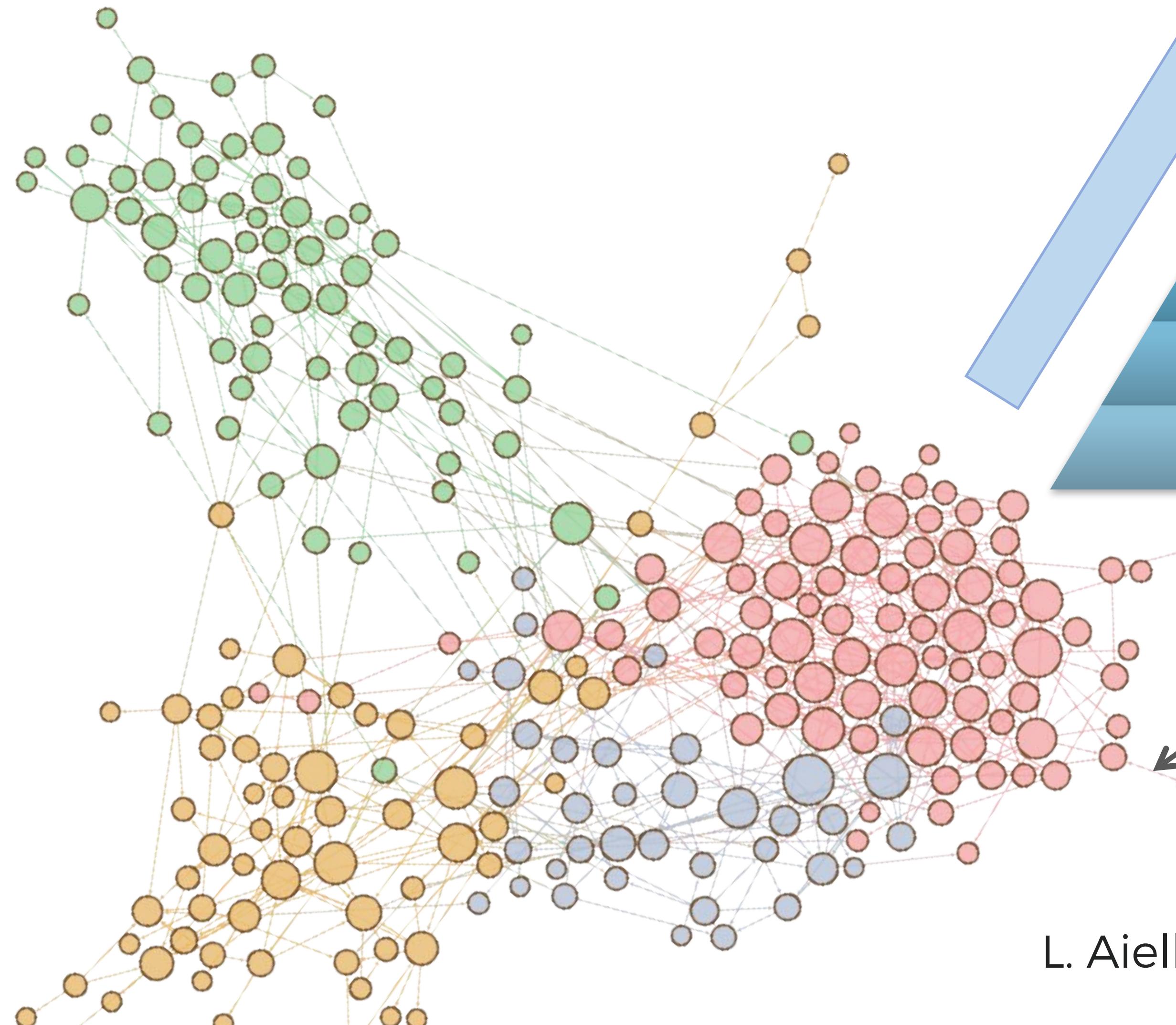


Social Exchange Theory

How we can model conversations?



Modeling social networks



Messages Are Communication Acts That
Define the Type of Social Relationship

+PRAGMATICS
(BEYOND SAYING)

Anobii.com dataset

 **Aliceassassina**
Female, 34. Roma, Italy

[Follow](#) [Message](#)

Friends [more](#)

-  [La Loulu](#)
-  [Tripwood76](#)
-  [Nerorossobia...](#)

Neighbors [more](#)

-  [Zoro](#)
-  [maldido duen...](#)
-  [ginocchiaapu...](#)

Shoutbox [See all](#)

Leave a comment

Morgana1981 "16 Novembre 2013 è uscito "Le Lenzuola nere di Khloe", Oct 28, 2013

CIAO!
Questo è il booktrailer del mio nuovo libro: "Le Lenzuola Nere di Khloe". Dura 2 minutini scarsi!
Buona Visione! :)

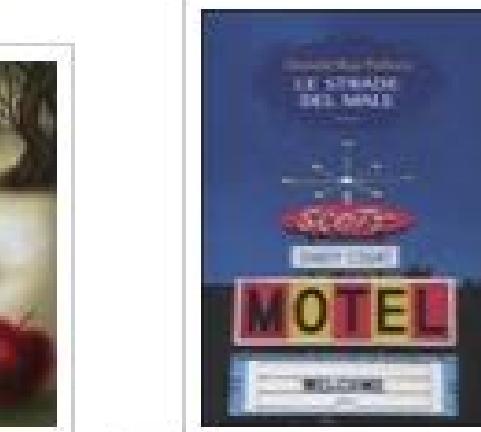
Books

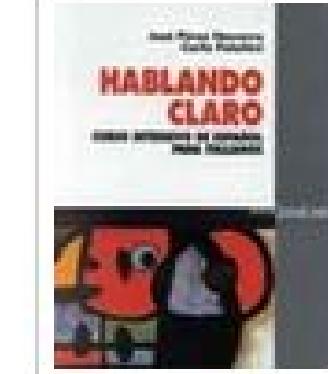
All books ▾ Rating (Highest)

 **L'Erbario delle Fate** (85)
By Benjamin Lacombe, Sébastien Perez
Finished on Nov 28, 2013  

 **La profezia dell'armadillo** (2931)
By Zerocalcare
Finished on Jan 6, 2013  

 **Biancaneve** (171)
By Jacob Grimm, Wilhelm Grimm
Finished on Feb 11, 2012  

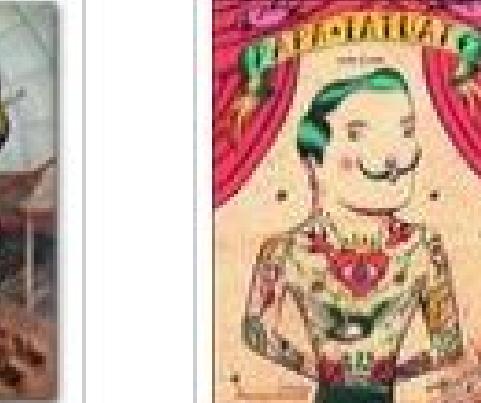
 **Le strade del male** (107)
By Donald Ray Pollock
Finished on Aug 17, 2012  

 **Hablando claro** (14)
Curso intensivo de Espanol para italianos
By Carla Polettini, José Pérez Navarro
Finished in 1999  

 **La ballata di Trenchmouth** (50)
Libro pop-up

 **C'era una volta...** (113)
Reference

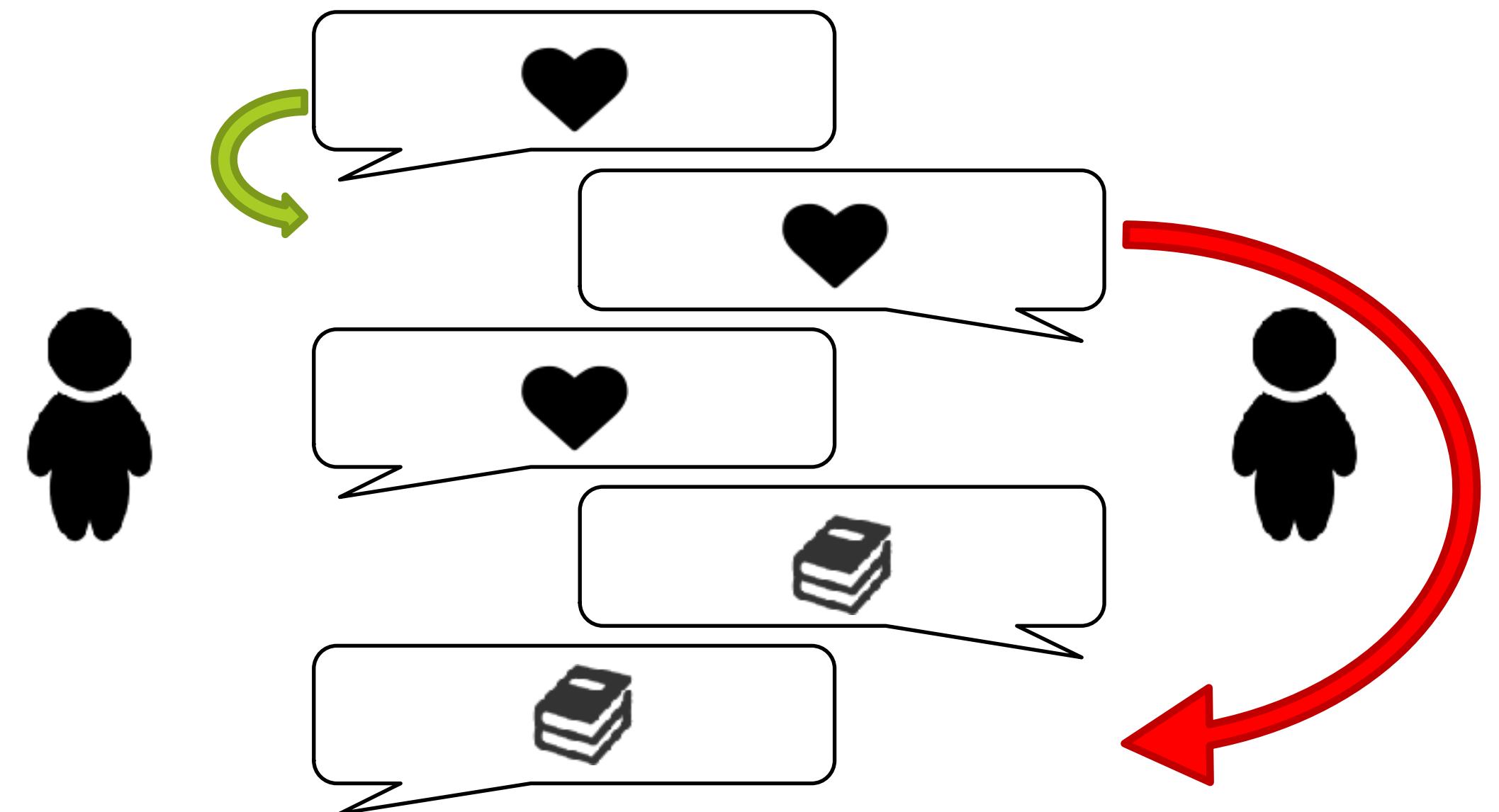
 **Circus Book 1870 - 1950** (4)
Reference

 **Papà tatuato** (67)
By Daniel

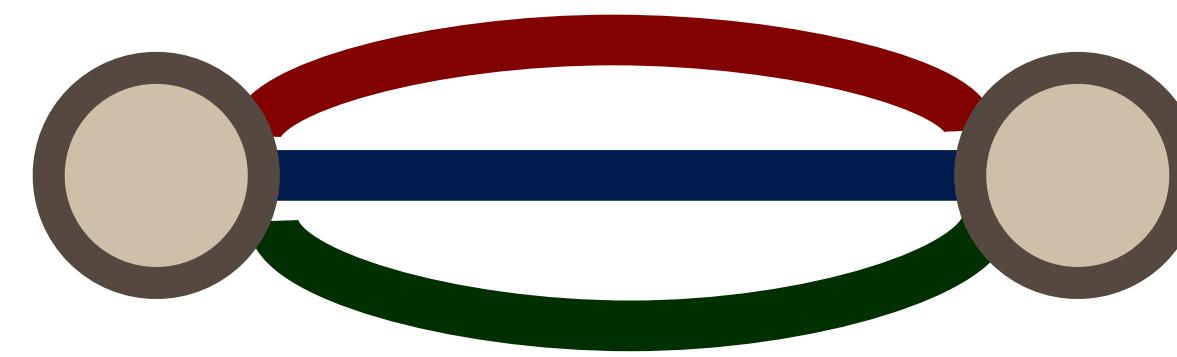
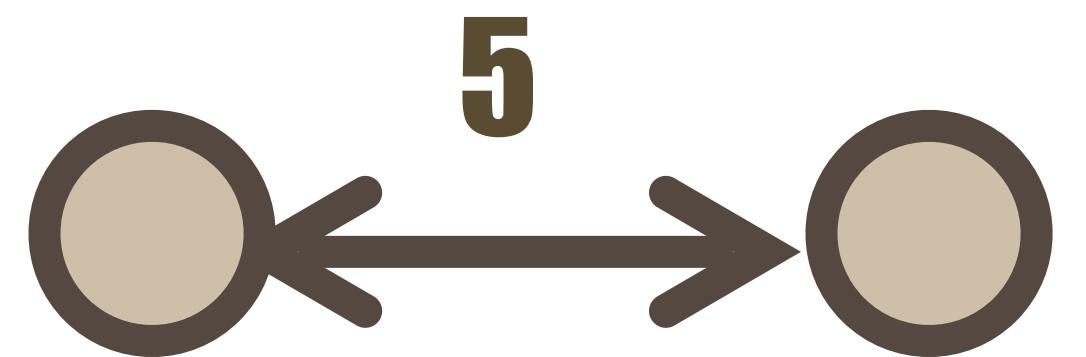
 **Venere privata** (2507)
By Giorgio

Blau's Social Exchange Theory

Repeated set of **exchanges** of **different types of non-material resources** transacted in an interpersonal situation, such as **knowledge, social support or manifestation of approval** [1]



[1] P. Blau. Exchange and Power in Social Life. Transaction Publishers, 1964.

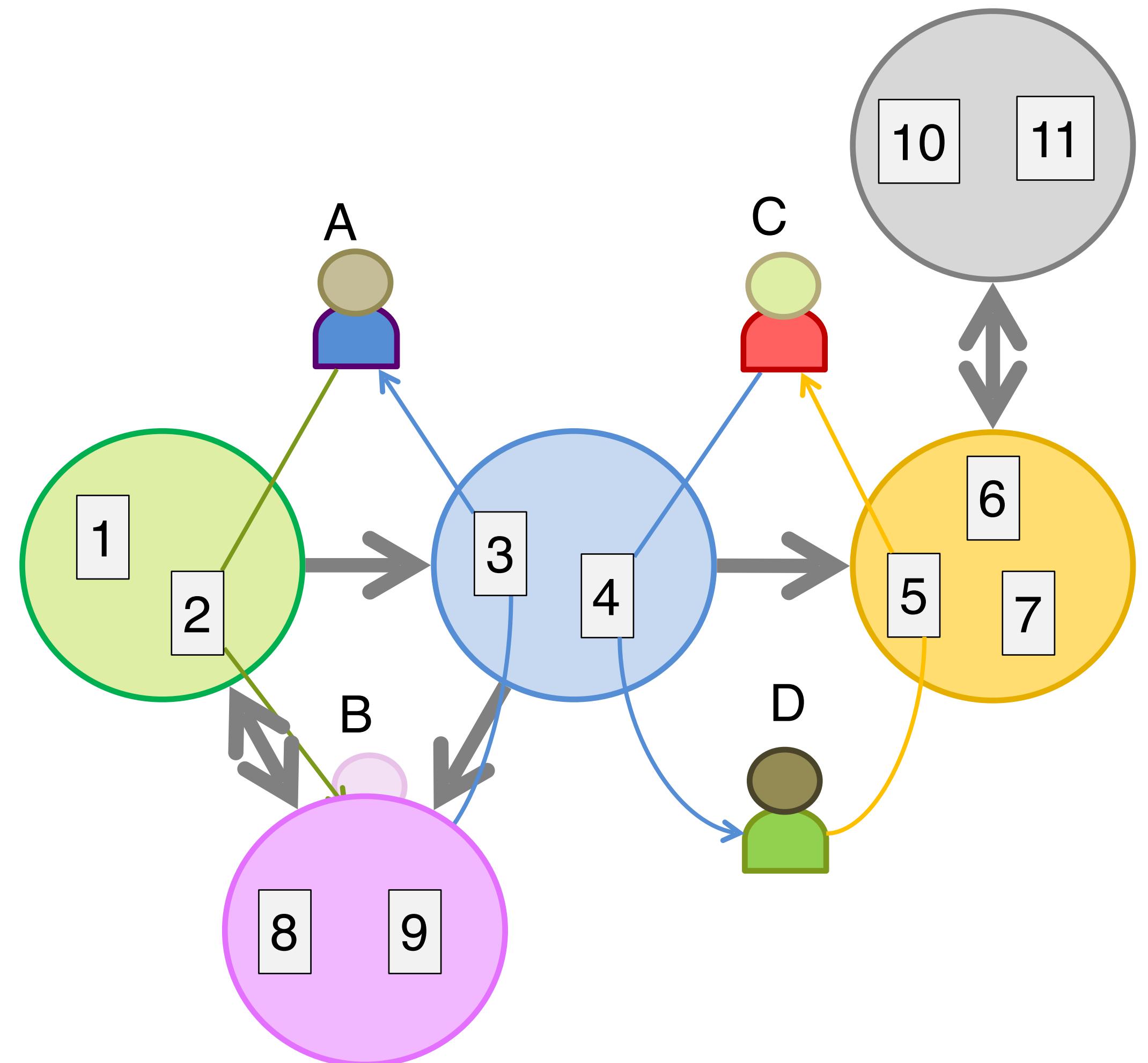


DECOMPOSITION OF A SOCIAL TIE INTO
the resources it conveys

How?

Input: directed comm. multigraph, arcs labeled with time and text
Output: (probabilistic) assignment message -> resource

- **Preprocessing**
 - Stopwords, stemming
- **Message bucketing**
 - NMF, LDA, ...
- **Transition graph**
 - Buckets as nodes transitions as edges
- **Intuition:** conversations tend to stick to the same resource (“You’re very good at it” -> “You are pretty good as well”)
- **Resource extraction**
 - Community detection on transition graph



Emerging resource types

Knowledge exchange

Technical knowledge of a domain (stackoverflow)

Request for knowledge

“I read a very good review of that book”

Status exchange

Expression of admiration or esteem

Recognition of the partner's higher status

“You are very smart!”

Social support

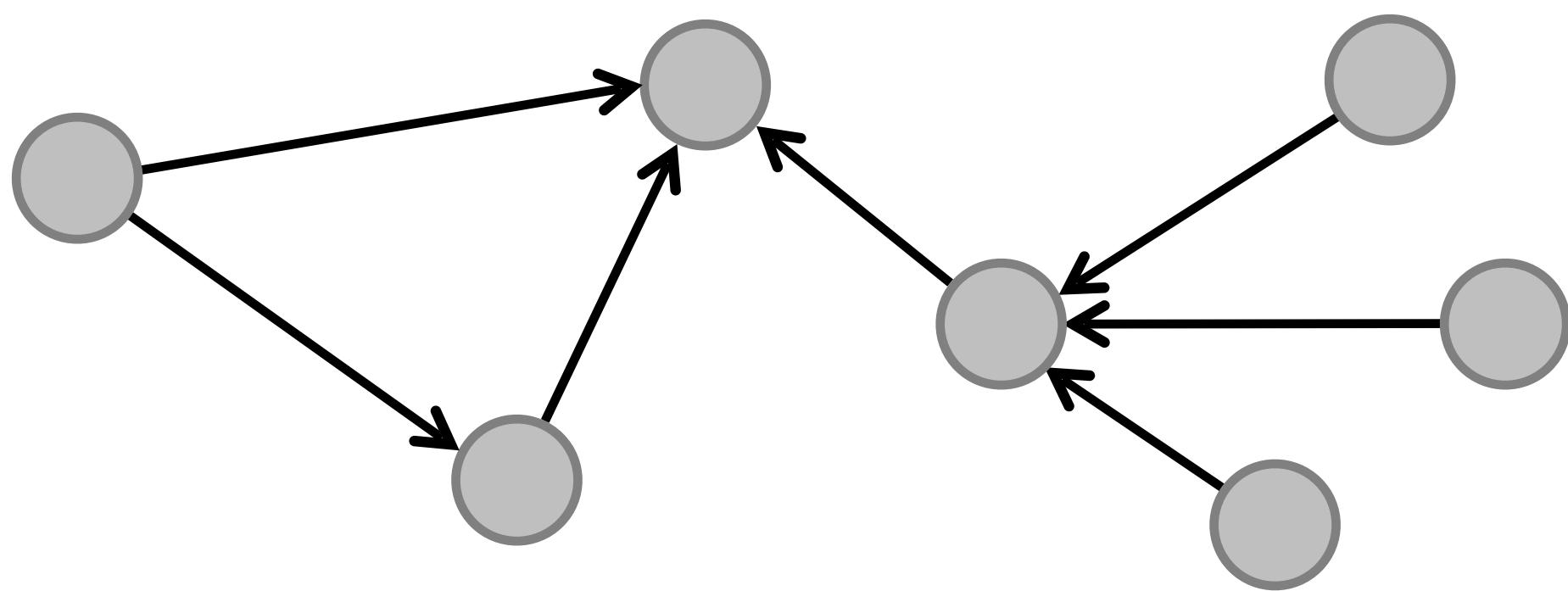
Emotional valuation

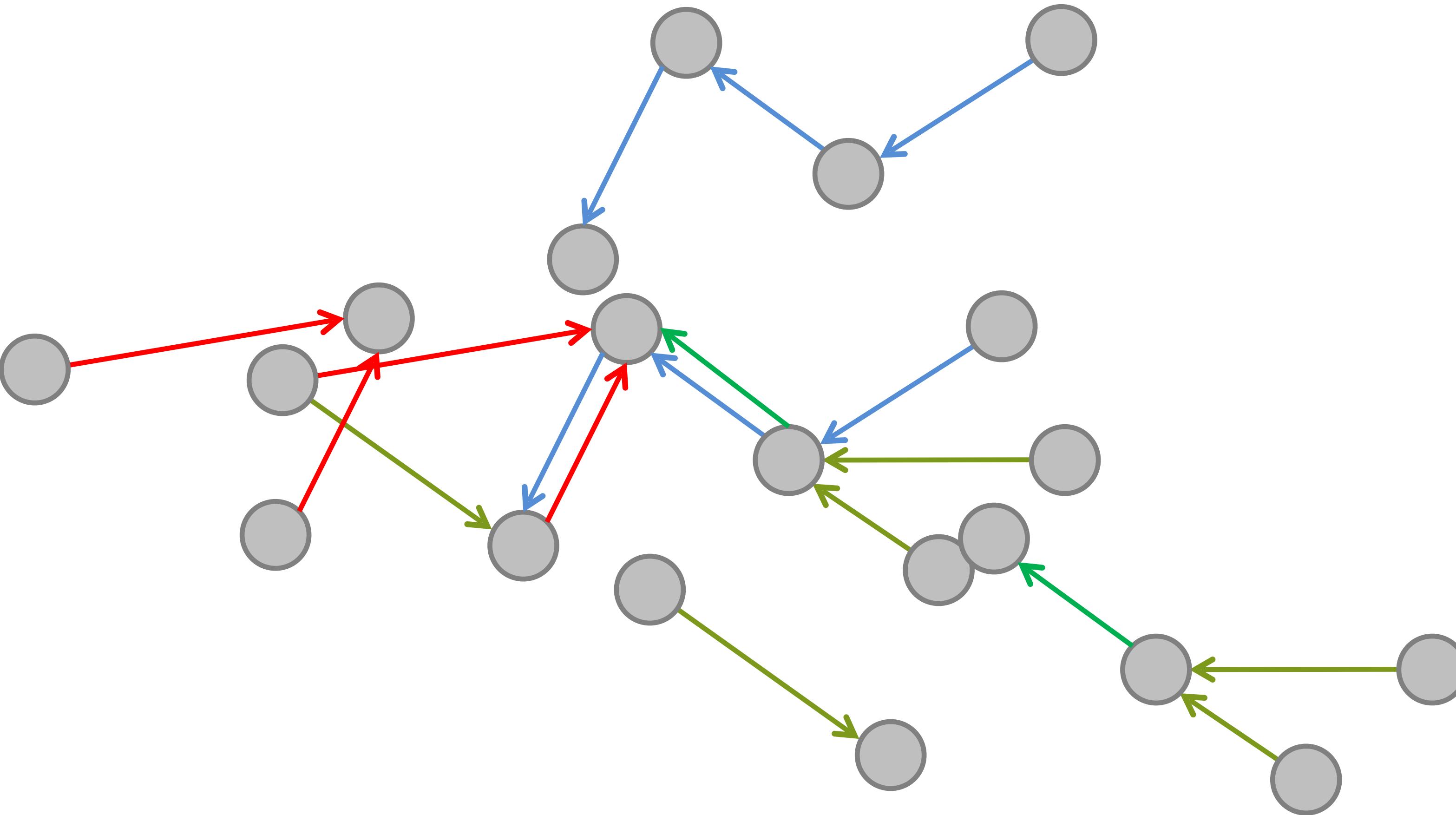
Everyday minute exchanges

“Hope your dad is feeling better now”

	Most representative tokens
<i>Status</i>	neighbor · <u>library</u> · return · read · bye thanks · much · hi · good read · thanks · <u>interest</u> · add · like · <u>congratulations</u> · visit
<i>Support</i>	good · morning · <u>wish</u> · <u>dear</u> · pass · good day · many · <u>greeting</u> · well · friend · soon · here · done · year · week · kiss
<i>Knowledge</i>	<u>book</u> · read · arrive · see · <u>know</u> · what · now · when · first · end · think · then · <u>advice</u> · only · again · have to · anobii · why

80% of messages are correctly assigned (human coders)

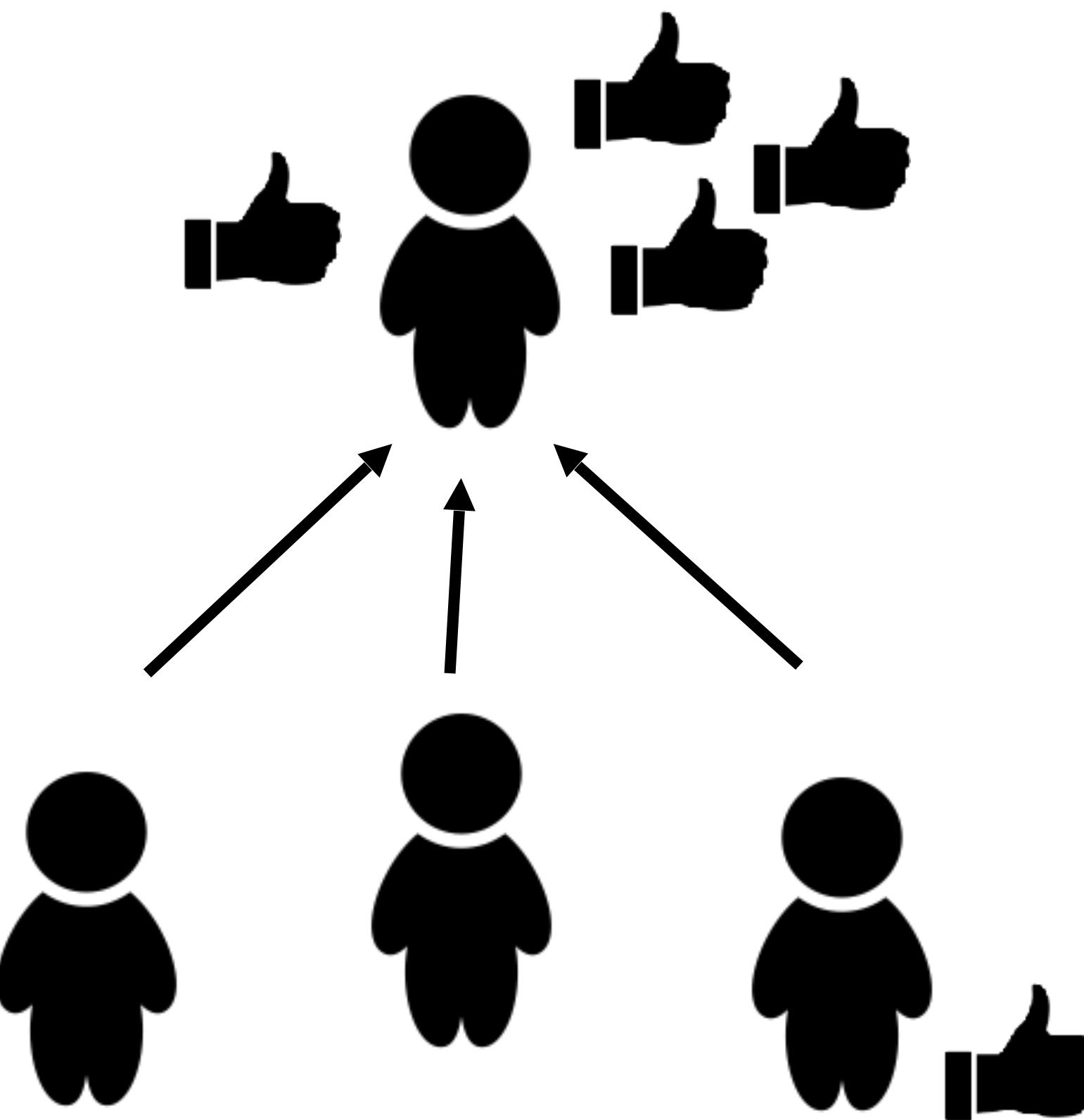
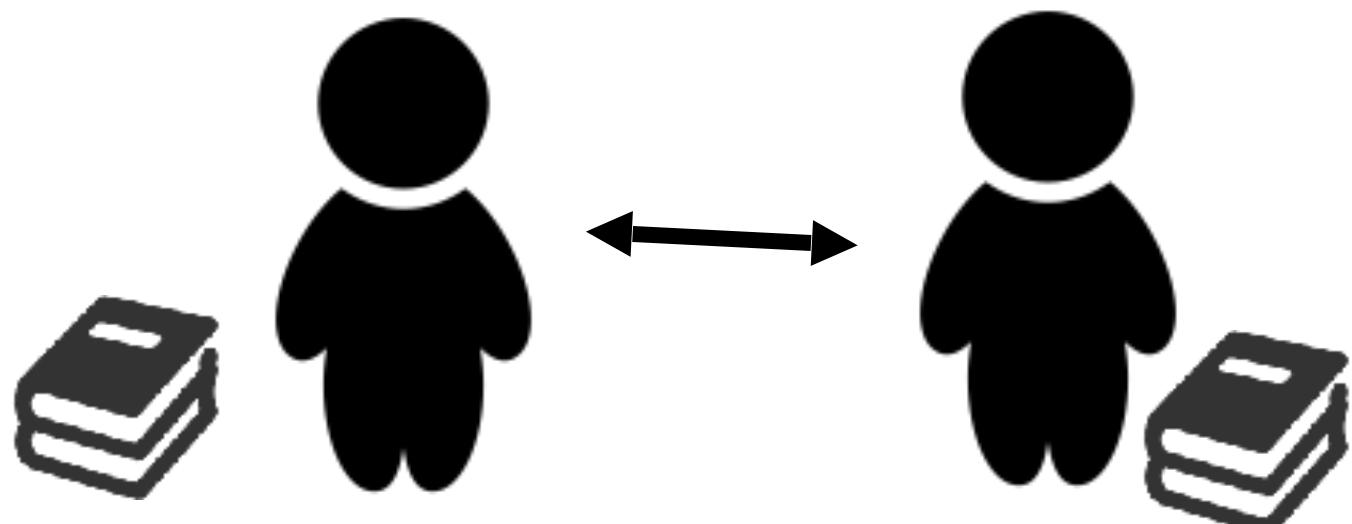
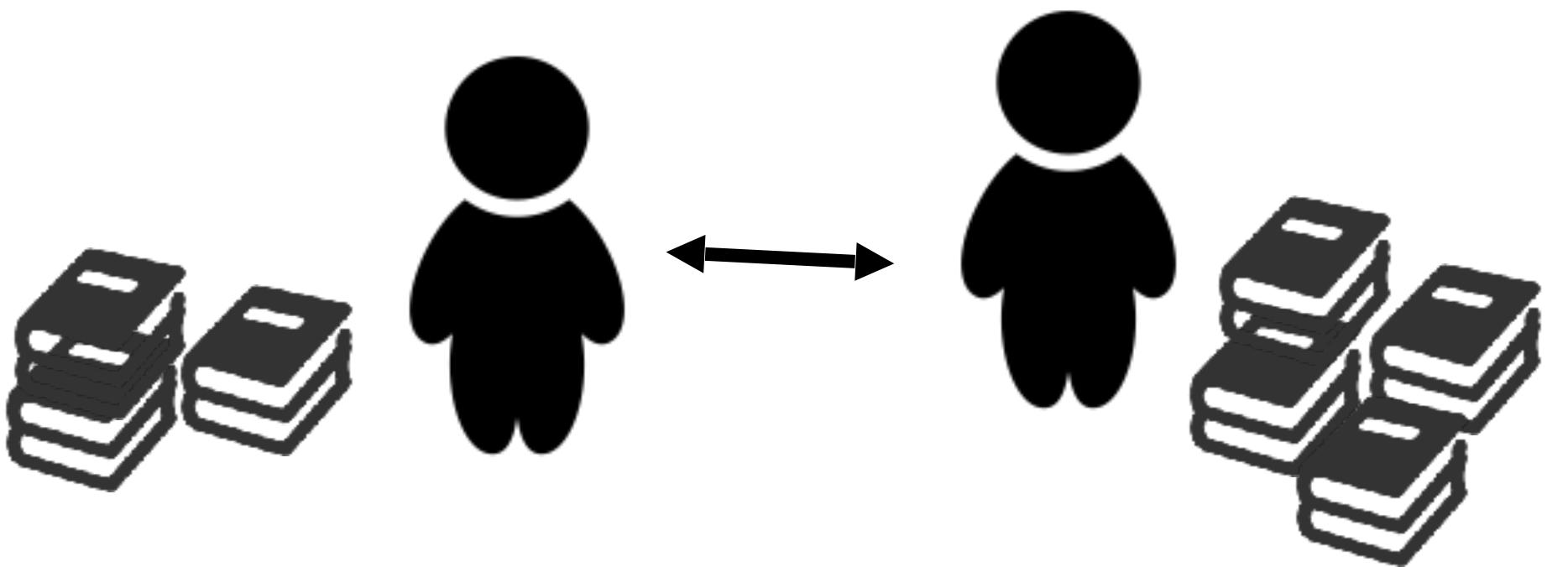




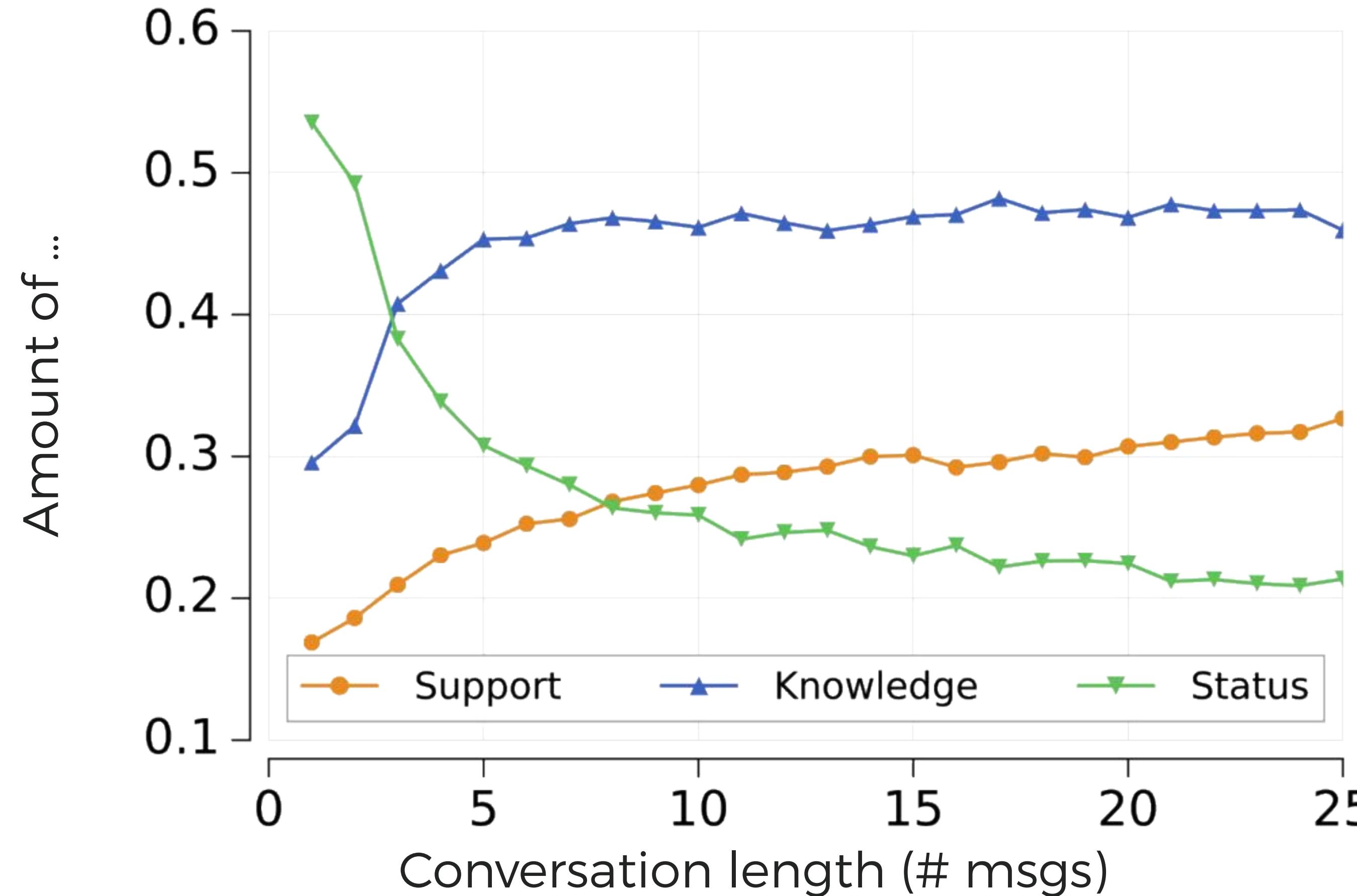
Knowledge net

VS.

Status net



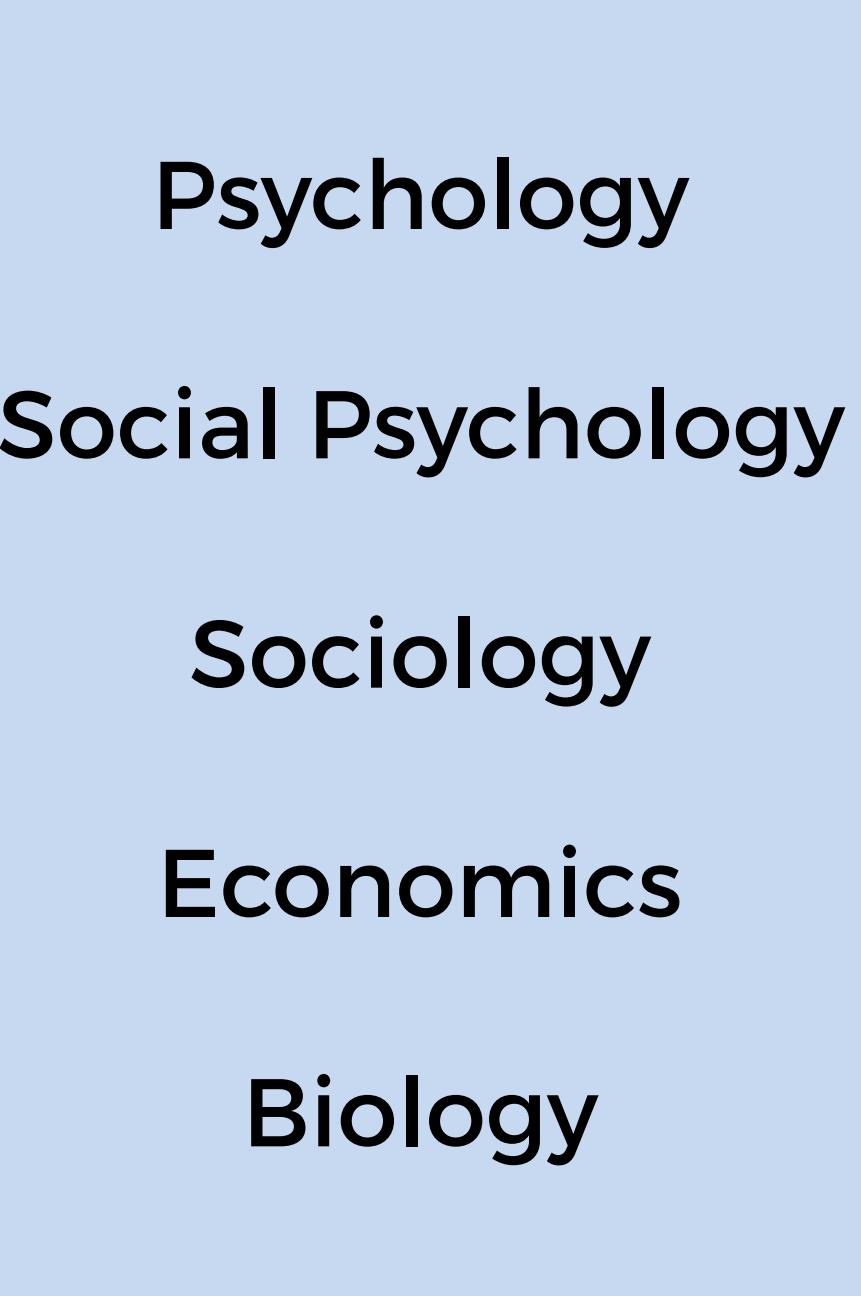
Resources over time





Literature review

Literature review



Psychology
Social Psychology
Sociology
Economics
Biology

~1k papers surveyed

Similarity

shared interests, hobbies, motivations, outlooks on life

Trust

high degree of dependability, and reliability

Romance

physical, emotional, and sexual attraction

Social support

high degree of emotional warmth, closeness

Identity

forms of meaningful common group memberships

Status

admiration, respect, prestige

Knowledge transfer

exchange of ideas, learning, and mentorship

Power

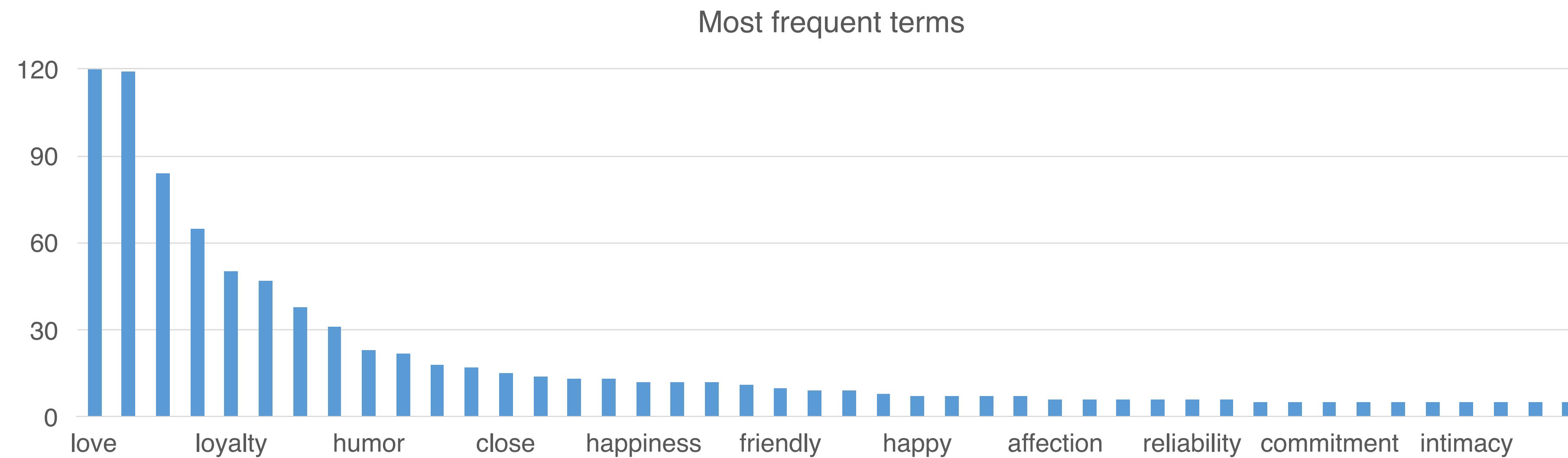
having the "upper hand," more resources, imbalance

A black and white photograph of a large lecture hall. The perspective is from an upper level, looking down at rows of students seated in tiered seating. The students are facing a professor who is standing at a podium on a stage. The stage has a large screen and some equipment. The overall atmosphere is that of a traditional university lecture.

crowdsourcing

Word-collection on MTurk

“Write the words the [best describe | matter most to you] in a social relationship”



362 Unique words, 151 non-redundant
200 workers. 45% women, 81% White from Canada, UK, US

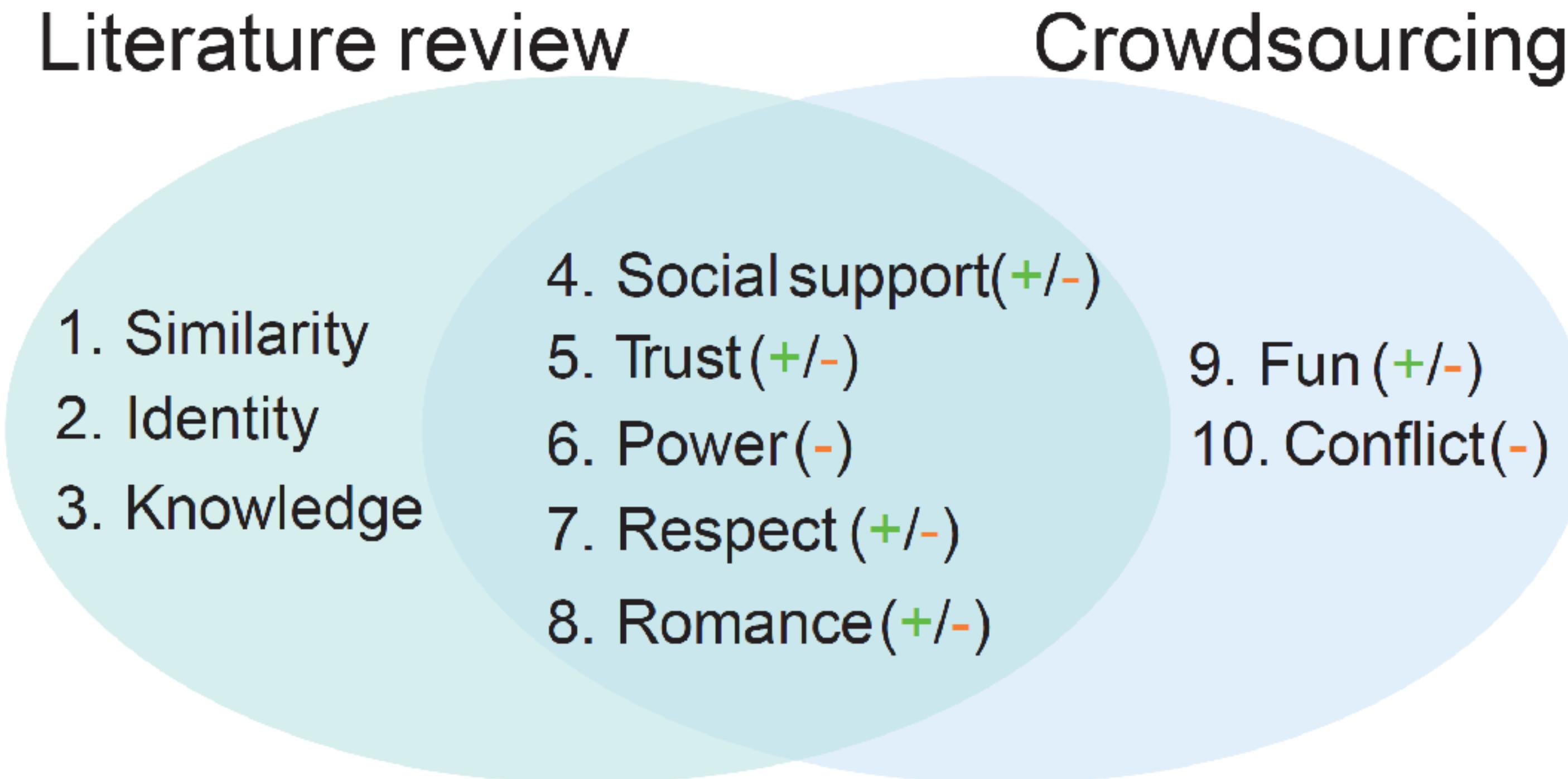
Word-scoring on MTurk

Think about the relationships in your life, in general (e.g. friends, family, coworkers, romantic partners). Now look at the words below. How well do the words below describe those relationships?

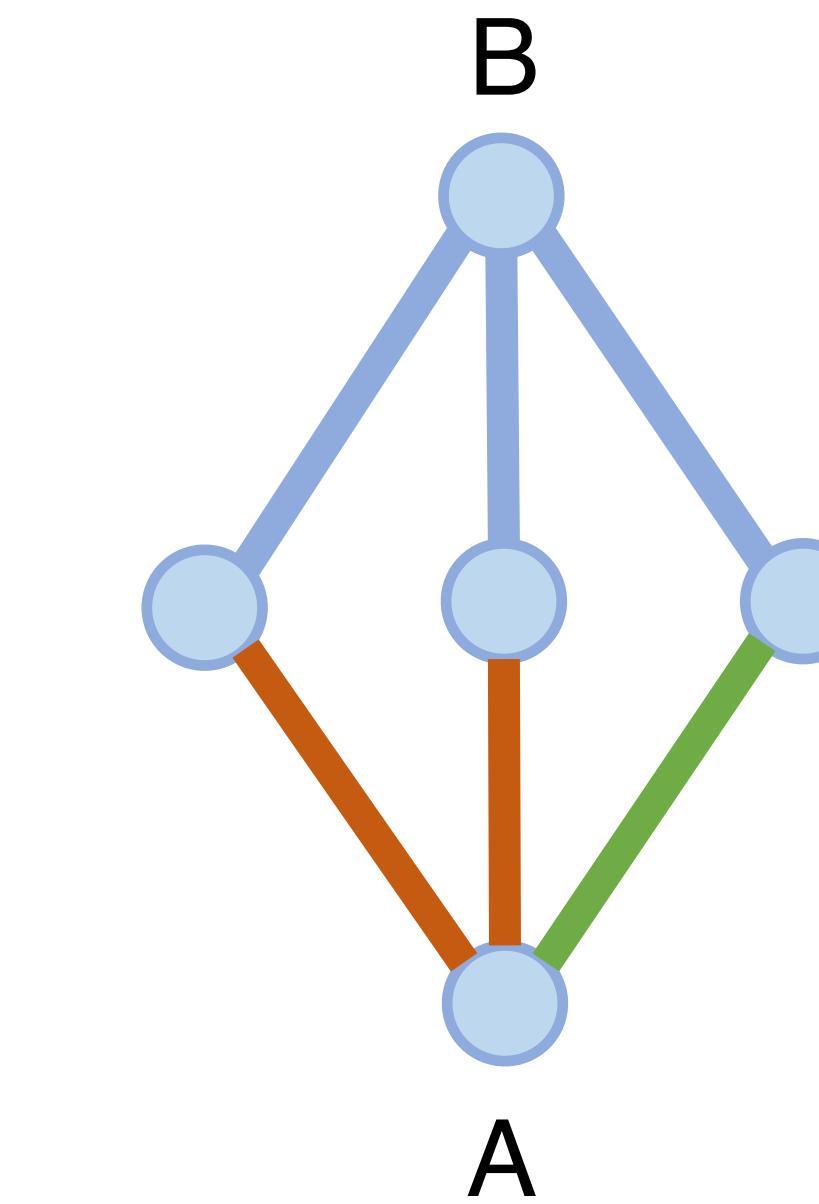
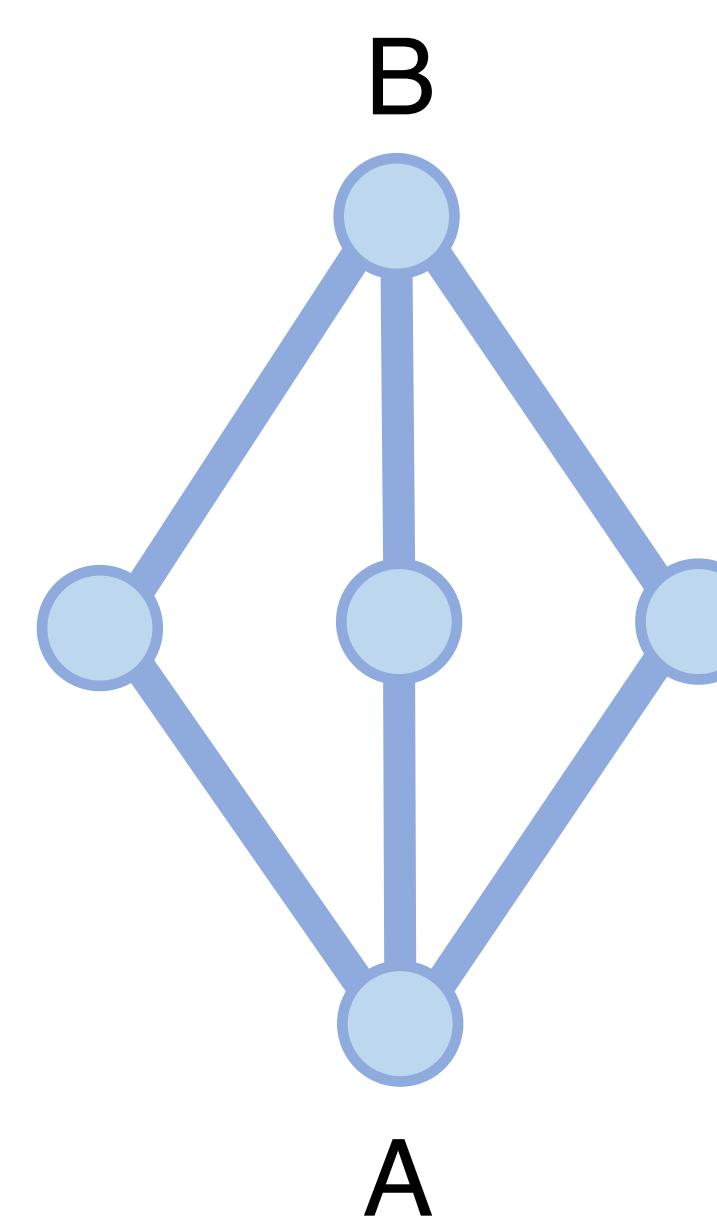
	1 Does not describe well at all	2	3 Describes fairly well	4	5 Describes very well
abusive	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
accountability	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
active	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
admiration	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
advantageous	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
adventurous	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
affection	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
amazing	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
amicable	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
attraction	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Each word is represented by a 200-dimensional vector of scores (1 entry per worker)

Word clustering



Link prediction



Simple word matching



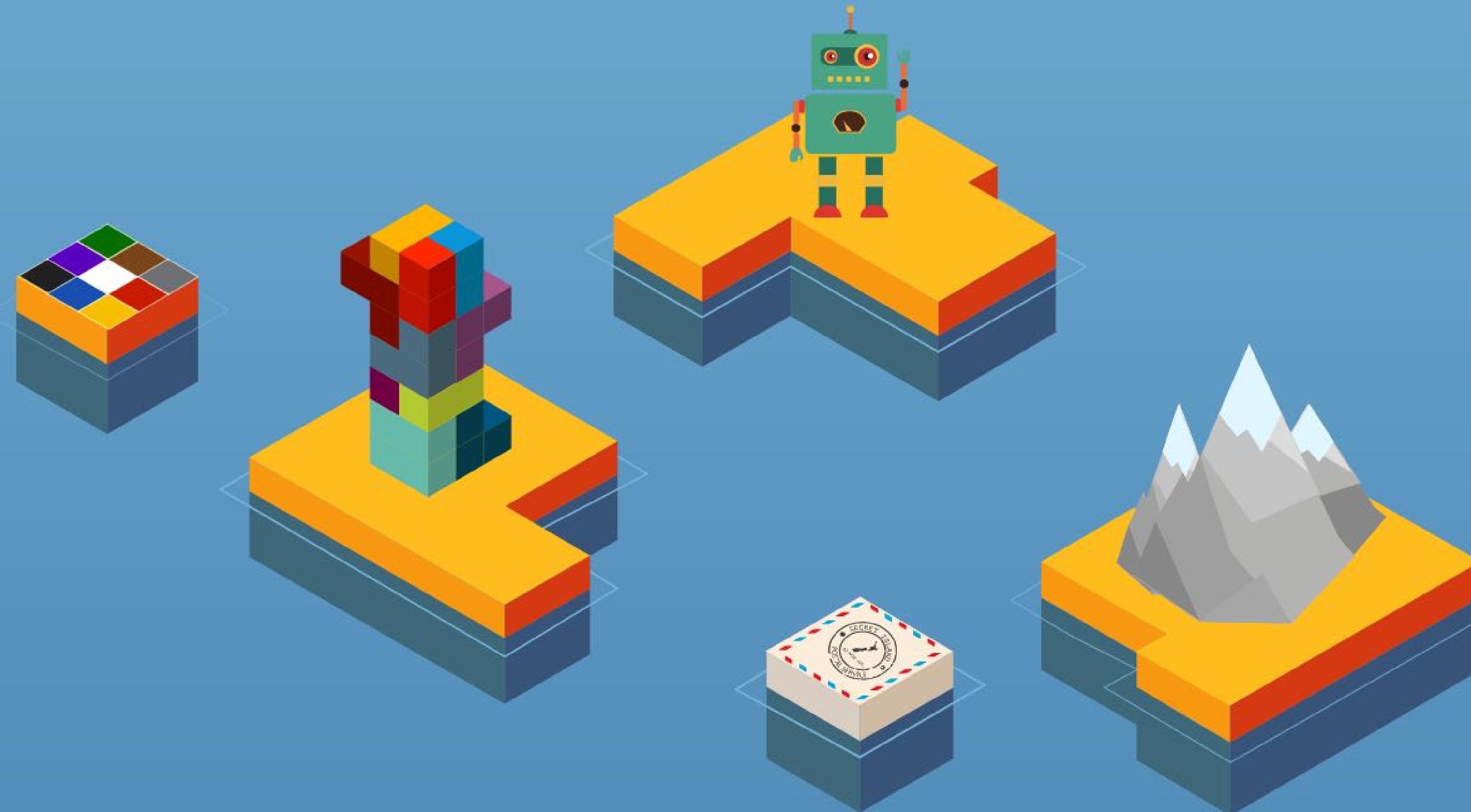
$$TO(u, v) = \frac{\Gamma_{out}(u) \cap \Gamma_{in}(v)}{|\Gamma_{out}(u)|}$$

Features	Precision	AUC
Triangle overlap	0.749	0.755
Relationship dimensions	0.800 (+7%)	0.803 (+6%)
All	0.783 (+5%)	0.841 (+11%)



f login

Click on an island and discover who you are



Click on an island and discover who you are



Isle of Ties

You likely know many people, but what are the relationships in your life really about? Play this game to find out what you really value in your relationships and how you compares to others.

 Connect to play

I'd rather play another game.

Pick three blocks describing your relation with *Alessandra Sala*

If you can think about one thing only, then pick the same block multiple times.

- Similarity
- Trust
- Romance
- Social support
- Identity
- Respect
- Knowledge transfer
- Power
- Fun
- Conflict



5/10



Undo

Skip

or restart



Isle of Ties

You likely know many people, but what are the relationships in your life really about? Play this game to find out what you really value in your relationships and how you compares to others.

Questions?

:::::::



@rschifan



schifane@di.unito.it



<http://www.di.unito.it/~schifane>