



---

# Social and Technological Network Data Analytics

## Lecture 1: Networks, Random Graphs and Metrics

Prof Cecilia Mascolo

# About Me



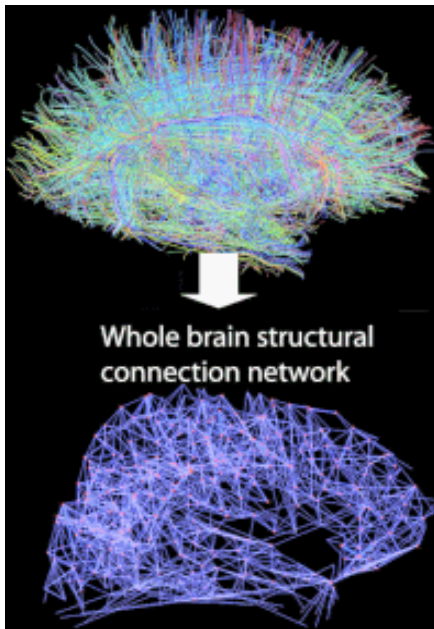
- Professor of Mobile Systems
  - NetOS Research Group
- Research on Mobile, Social and Sensor Systems
- More specifically,
  - Human Mobility and Social Network modelling
  - Geo-social recommendation systems
  - Mobile Sensor Systems and Networks



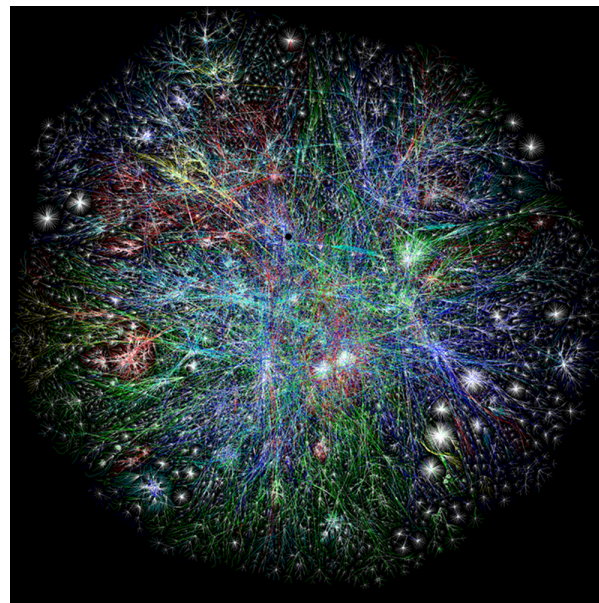
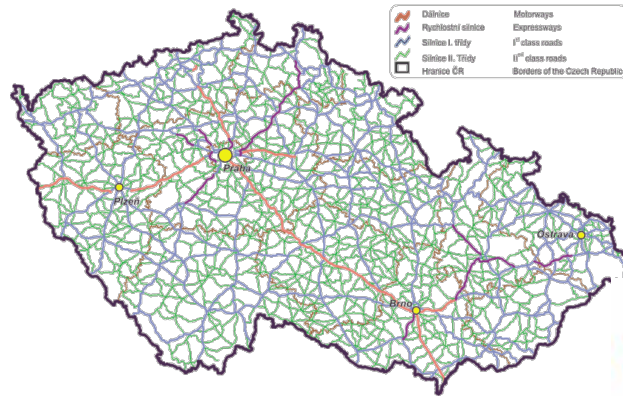
twitter



# Networks are Everywhere



Whole brain structural connection network



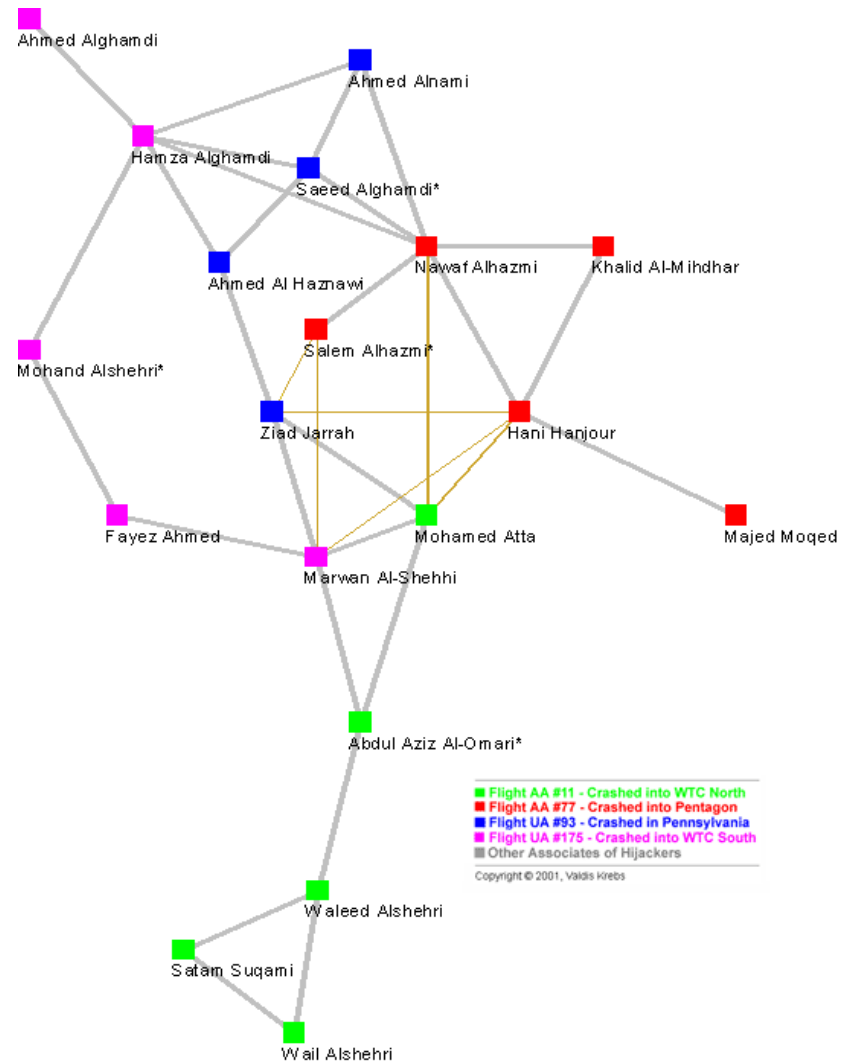
# Facebook Friendship Network



# Terrorist Network

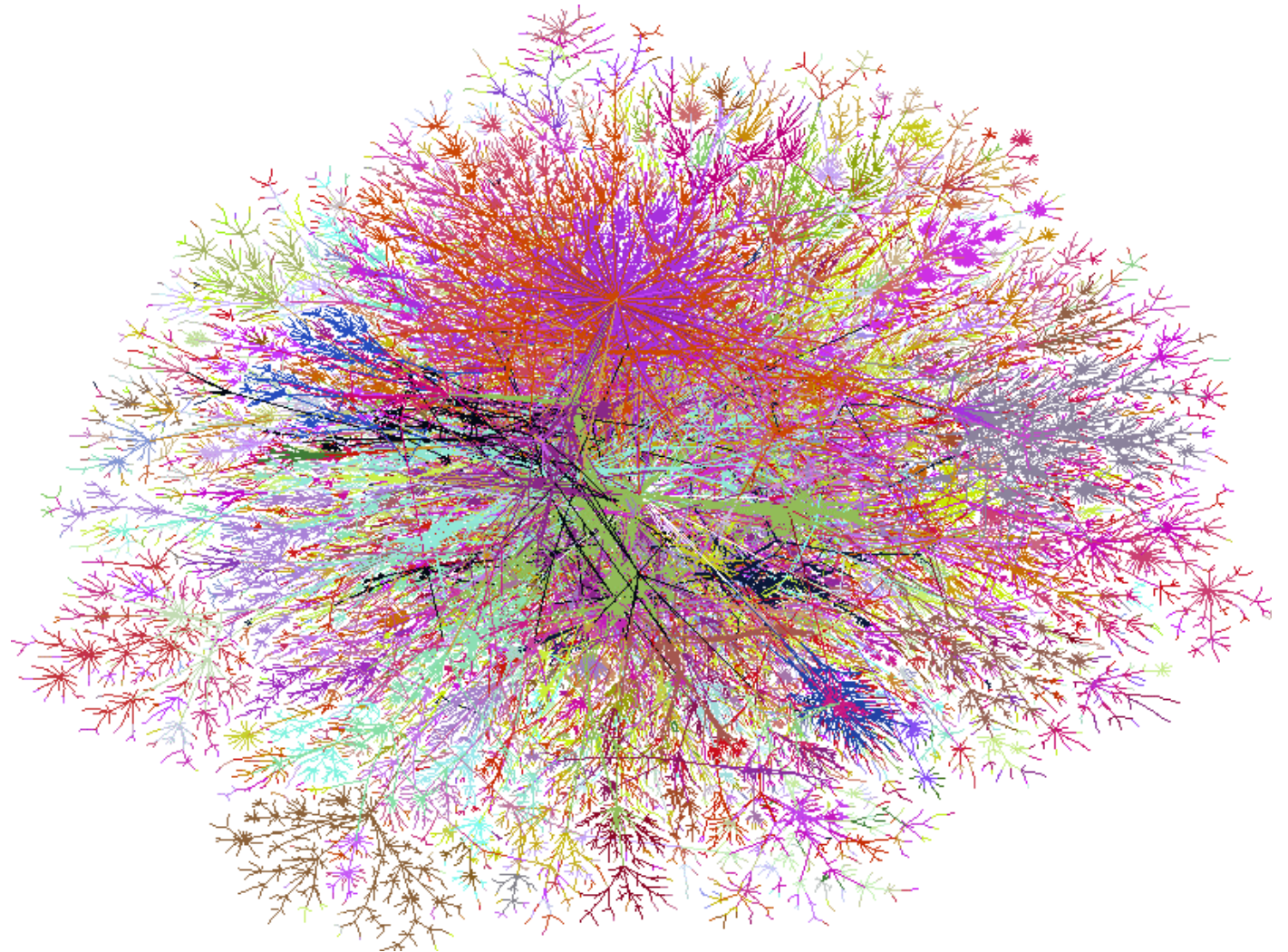


“Six degrees of Mohammed Atta”

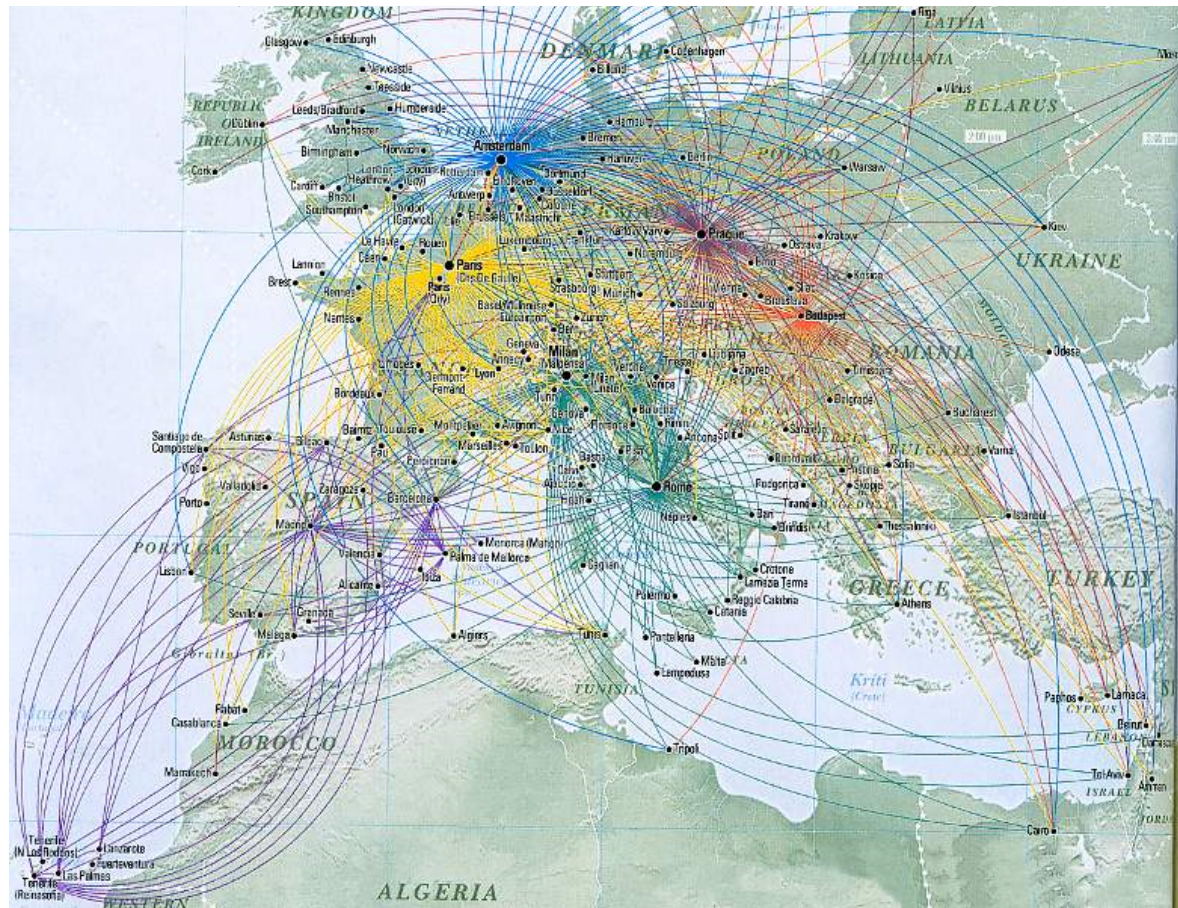


# The Internet

---



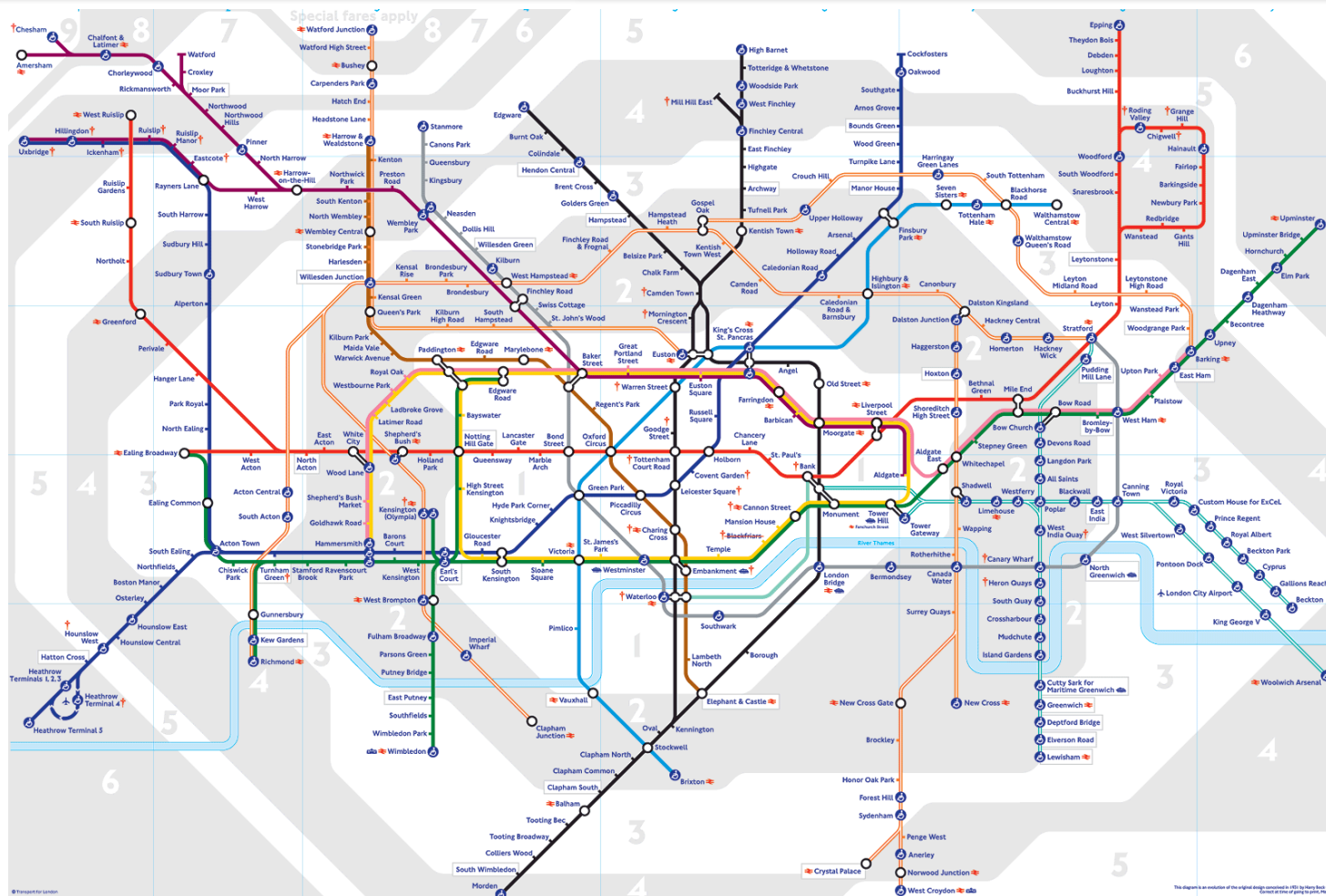
# Airline Network



UNIVERSITY OF  
CAMBRIDGE

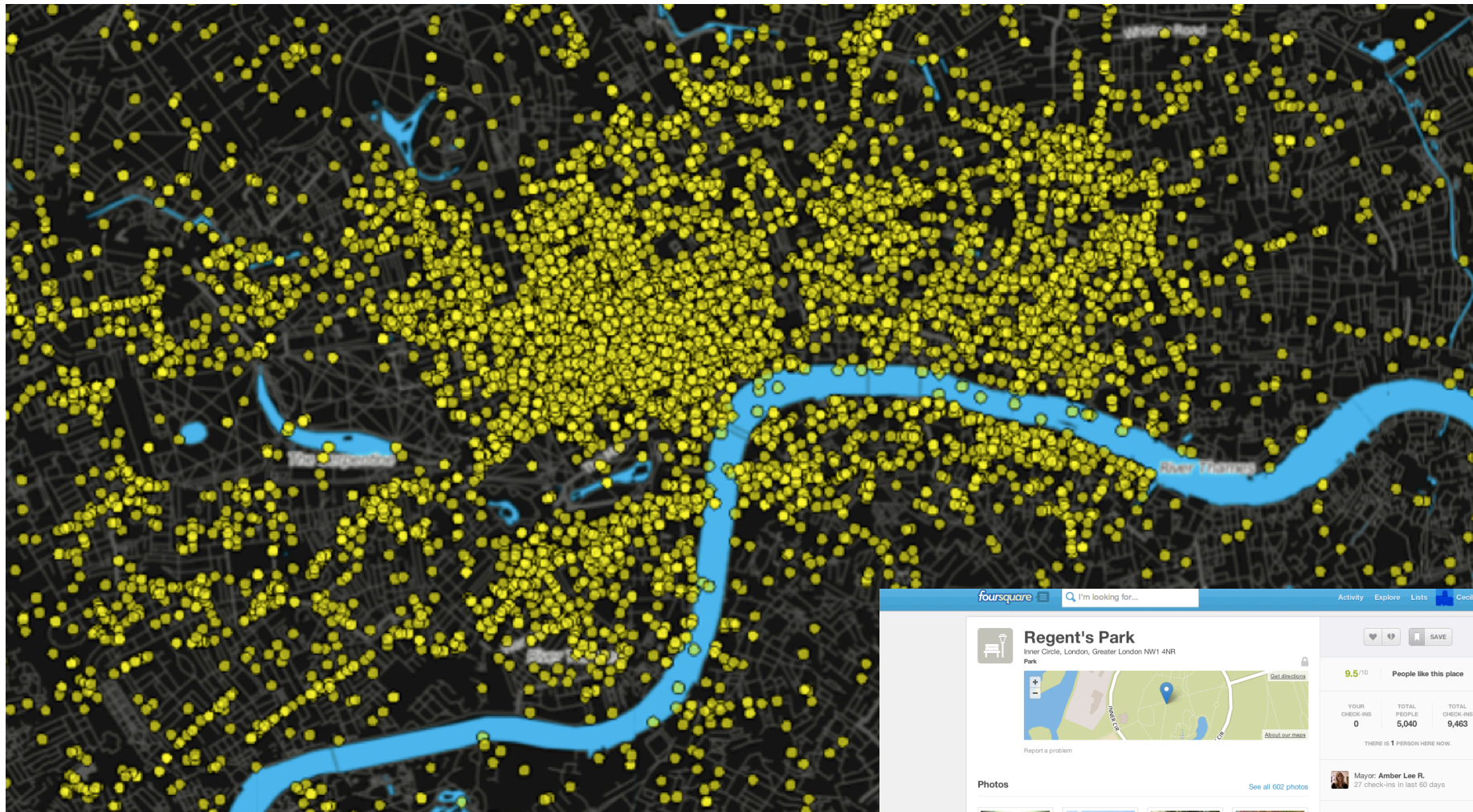
Source: Northwest Airlines WorldTraveler Magazine

# Railway/Metro Network





# Geo-social Networks



foursquare I'm looking for... Activity Explore Lists Cecils

### Regent's Park

Inner Circle, London, Greater London NW1 4NR

9.5/10 People like this place

YOUR CHECK-INS	TOTAL PEOPLE	TOTAL CHECK-INS
0	5,040	9,463

THERE IS 1 PERSON HERE NOW.

Mayor: Amber Lee R.  
27 check-ins in last 60 days

1 friend has been here

Similar places  
Primrose Hill, St James's Park, Hampstead Heath, Soho Square, Green Park

Explore Nearby  
Restaurants, Nightlife, Shopping, Top Picks

# What Kind of Networks?

---



- Who talks to whom?
- Who is friend with whom?
- What leads to what?
- Who is a relative of whom?
- Who eats whom?
- Who sends messages to whom?



# In This Course

---

- We will study the models and metrics which allow us to understand these phenomena.
- We will show analysis over large datasets of real social and technological networks.



# List of Lectures

---

- Lecture 1: Networks and Random Graphs
- Lecture 2: Small World and Weak Ties
- Lecture 3: Centrality and Applications
- Lecture 4: Community Detection, Modularity, Overlapping Communities
- Lecture 5: Structure of the Web, Search and Power laws
- Lecture 6: Network Robustness and Applications
- Lecture 7-8: 2h Practical Tutorial
- Lecture 9: Cascade and Behaviour Influence
- Lecture 10: Epidemic Spreading and Examples
- Lecture 11: Influence and Epidemic Spreading Applications
- Lecture 12: Spatial Networks
- Lecture 13: Temporal Networks
- Lecture 14: External Presentation
- Lecture 15-16: Student Presentations

# Assessment



- All information on the course page:  
<http://www.cl.cam.ac.uk/teaching/1617/L109/materials.html>
- A report of at most **4000 words** which consists of analysis of an assigned dataset according to some indicated network measures: the results should be commented and justified (**70% of the final mark**). The report is due on **13<sup>th</sup> March (noon)**
- 8 minutes presentation of the findings of the assignment on **14th March**. The presentation is **worth 30%** of the final mark.

# Dataset Assignment

---



- Send me email immediately with a choice of dataset (**and a backup choice**) from the website [assignment is first come first served].
- At that point you will receive an email with a link to the dataset.
  - **Do this by 10<sup>th</sup> February.**



# In This Lecture

---

- We will introduce:
  - Networks/graphs
  - Basic network measures
  - Random Graphs
  - Examples

# A Network is a Graph

---



A **graph**  $G$  is a tuple  $(V,E)$  of a set of vertices  $V$  and edges  $E$ . An edge in  $E$  connects two vertices in  $V$ .

A **neighbour set**  $N(v)$  is the set of vertices adjacent to  $v$ :

$$N(v) = \{u \in V \mid u \neq v, (v,u) \in E\}$$

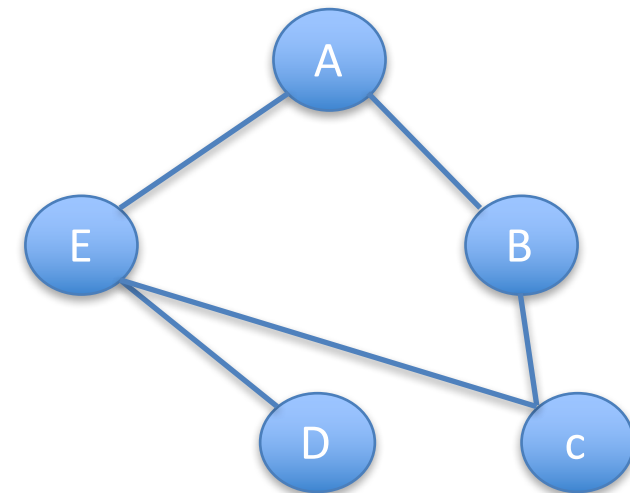




# Node Degree

- The **node degree** is the number of neighbours of a node
- E.g., Degree of A is 2

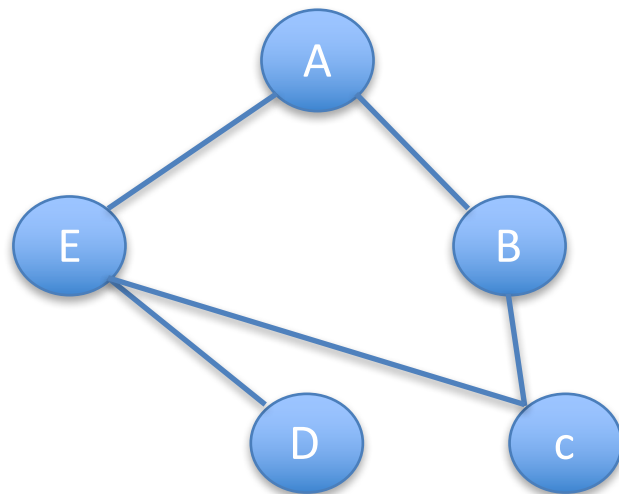
The study of the degree distribution of networks allows the classification of networks in different categories



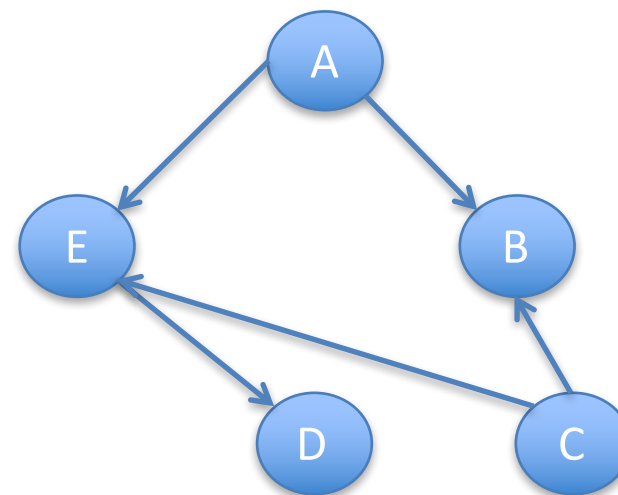


# Directed & Undirected Graphs

---



Undirected Graph



Directed Graph

Example of Undirected Graphs: Facebook, Co-presence

Examples of Directed: Twitter, Email, Phone Calls

# Paths and Cycles

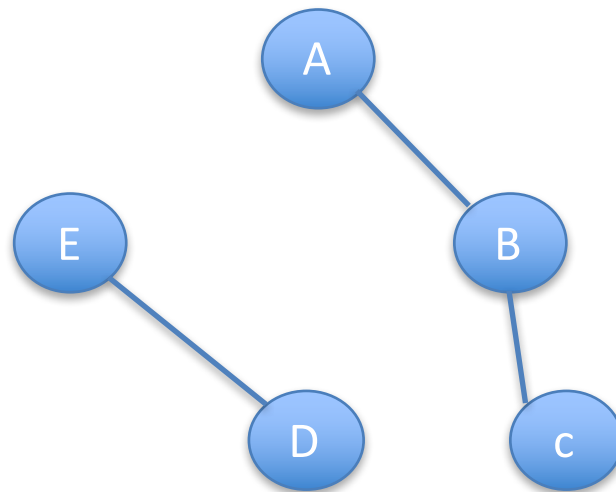


- A **path** is a sequence of nodes in which each pair of consecutive nodes is connected by an edge.
  - If a graph is directed the edge needs to be in the right direction.
  - E.g. A-E-D is a path in both previous graphs
- A **cycle** is a path where the start node is also the end node
  - E.g. E-A-B-C is a cycle in the undirected graph

# Connectivity



- A graph is **connected** if there is a path between *each pair* of nodes.
- Example of **disconnected** graph:



# Components

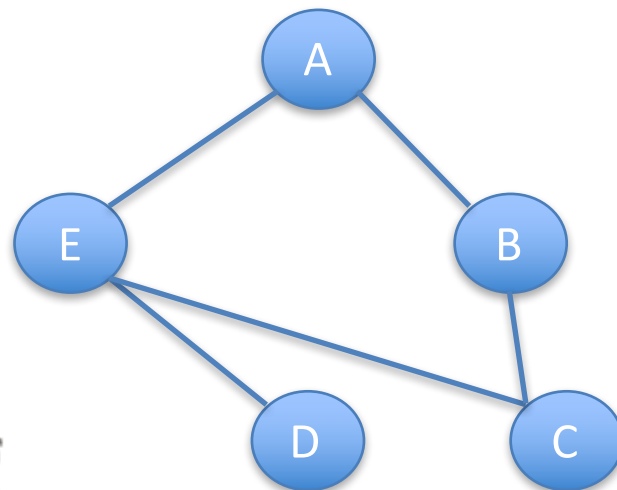


- A **connected component** of a graph is the subset of nodes for which each of them has a path to all others (and the subset is not part of a larger subset with this property).
  - Connected components: A-B-C and E-D
- A **giant component** is a connected component containing a significant fraction of nodes in the network.
  - Real networks often have one unique giant component.

# Path Length/Distance



- The **distance** ( $d$ ) between two nodes in a graph is the length of the shortest path linking the two graphs.
- The **diameter** of the graph is the maximum distance between any pair of its nodes.



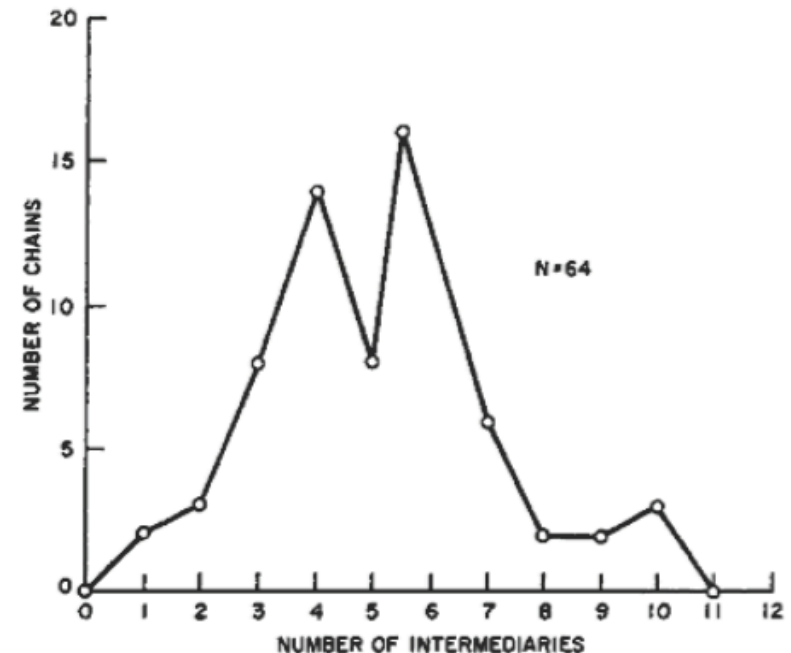
What's the diameter here?

# Small-world Phenomenon

## Milgram's Experiment



- Two random people are connected through only a few (6) intermediate acquaintances.
- Milgram's experiment (1967) shows the known "six degrees of separation":
  - Choose 300 people at random
  - Ask them to send a letter through friends to a stockbroker near Boston.
  - 64 successful chains.



# Recent Remake



- In 2003 the experiment was redone.
- Choice of 18 targets over the world. Choice of senders (60K) from a commercially obtained email list.
- Website to control the “email” contact from one participant to the next (two weeks to choose next hop)
- Verification of relationship by receiver (to avoid cheating by web search).

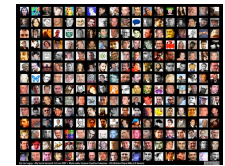


# Findings



- Use of “weak ties” and professional relationships
- Median of 5-7 steps
- “Network structure alone is not everything”
- Some different incentives had a high impact on completion rate of chains
  - If the target was in a prominent place (eg professor)

# Milgram's study on Facebook



facebook



Search for people, places and things



Cecilia Mascolo

Home



## Anatomy of Facebook

By Lars Backstrom on Tuesday, 22 November 2011 at 01:04

Think back to the last time you were in a crowded airport or bus terminal far from home. Did you consider that the person sitting next to you probably knew a friend of a friend of a friend of yours? In the 1960s, social psychologist Stanley Milgram's "small world experiment" famously tested the idea that any two people in the world are separated by only a small number of intermediate connections, arguably the first experimental study to reveal the surprising structure of social networks.

With the rise of modern computing, social networks are now being mapped in digital form, giving researchers the ability to study them on a much grander scale than ever before.



Notes by Facebook Data Science

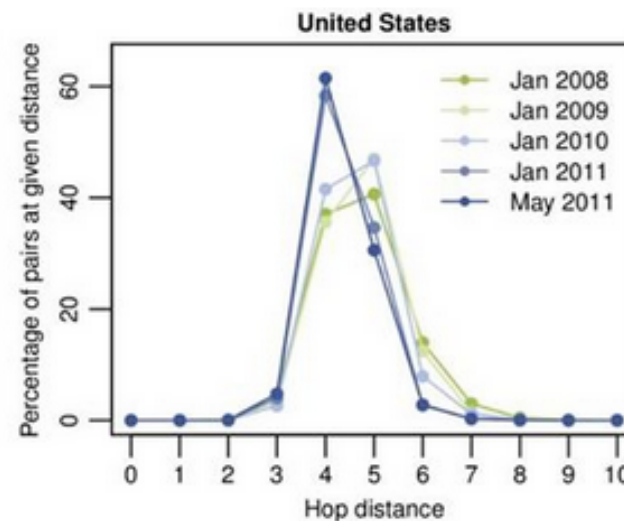
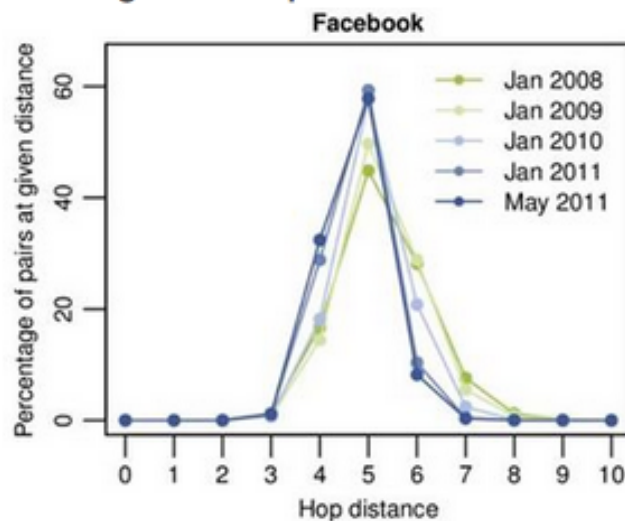
All Notes

Get notes via RSS

Embed Post

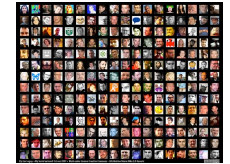
Report

### Four degrees of separation.

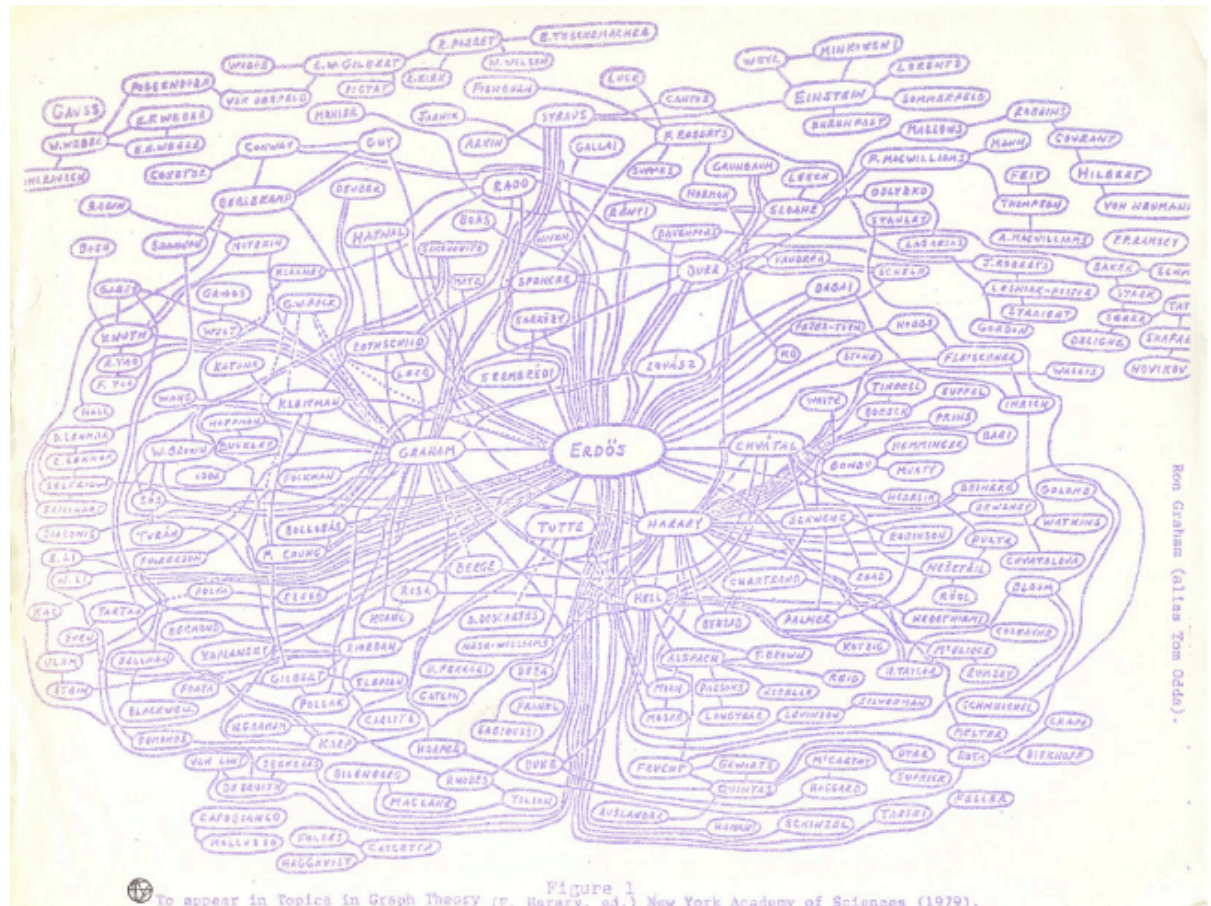


While 99.6% of all pairs of users are connected by paths with 5 degrees (6 hops), 92% are connected by only four degrees (5 hops)

# Erdos Number



- Erdos Number: distance from the mathematician (most people are 4-5 hops away) based on collaboration.



# Bacon Number



- A network of actors who costarred in a movie.
- Most actors are no more than  $\sim 3$  hops from Kevin Bacon.
- One very obscure movie was at distance 8.

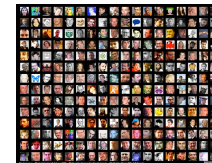


# Random Graphs



- First way to model these networks:
- **Erdos-Renyi Random Graph** [Erdos-Renyi '59]:

$G(n,p)$ : graph with  $n$  vertices where an edge exists with independent random probability  $0 < p < 1$  for each edge.



# Random Graph Model

- For each node  $n_1$ , an edge to node  $n_2$  exists with probability  $p$ .
- Degree distribution is **binomial**.
- The probability of a node to have degree  $k$ :

$$P(k_i = k) = C_{N-1}^k p^k (1-p)^{N-1-k}$$

- Where  $C_{N-1}^k = \binom{N-1}{k}$
- Expected Degree of a node:  $(N-1)p \approx Np$



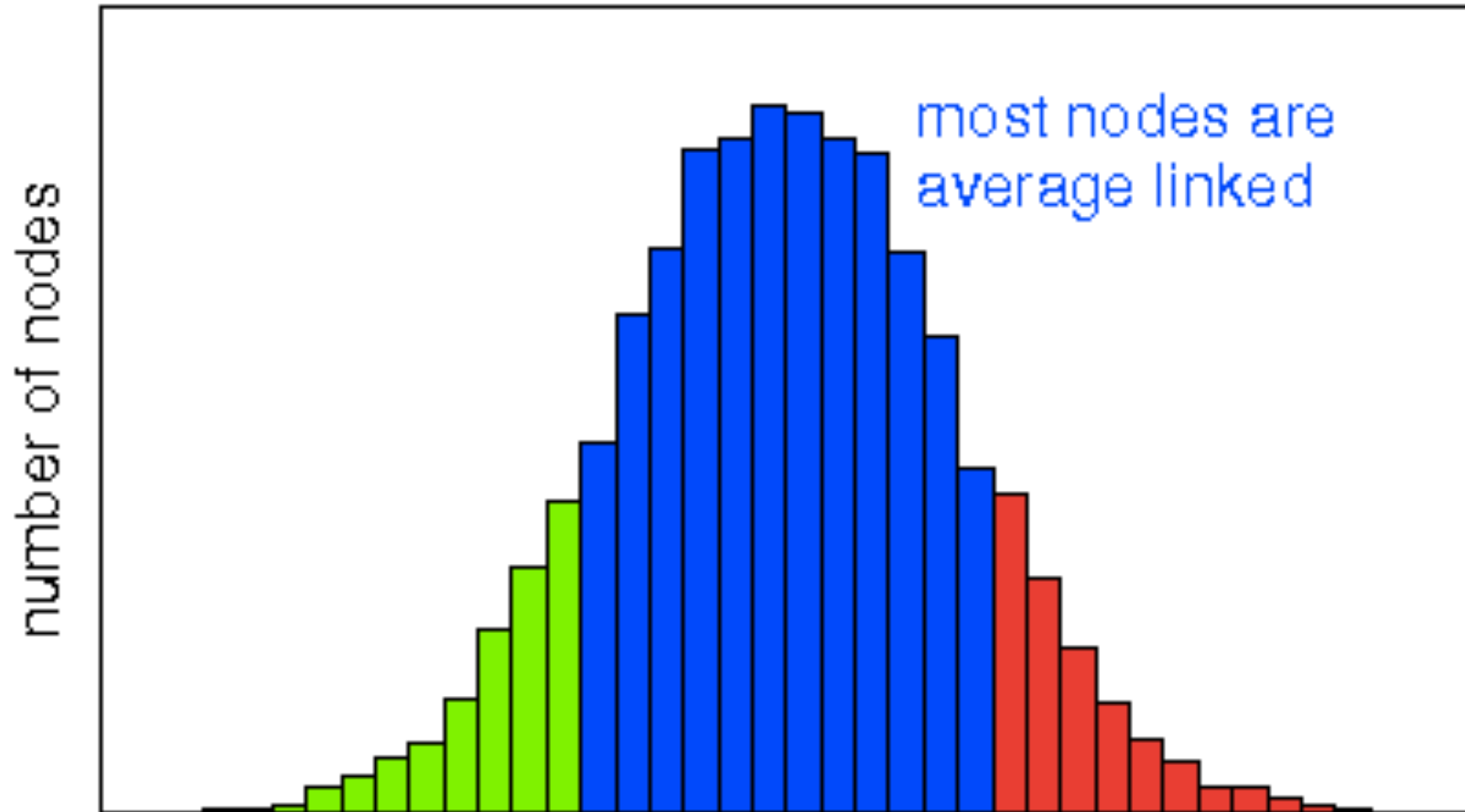
# Random Graphs Properties

---

- For large  $N$  this is approximated by the Poisson distribution with

$$P(k) \approx e^{-Np} \frac{(Np)^k}{k!} = e^{-\langle k \rangle} \frac{\langle k \rangle^k}{k!}$$

# Degree Distribution of Random Graphs





# Random Graph Diameter

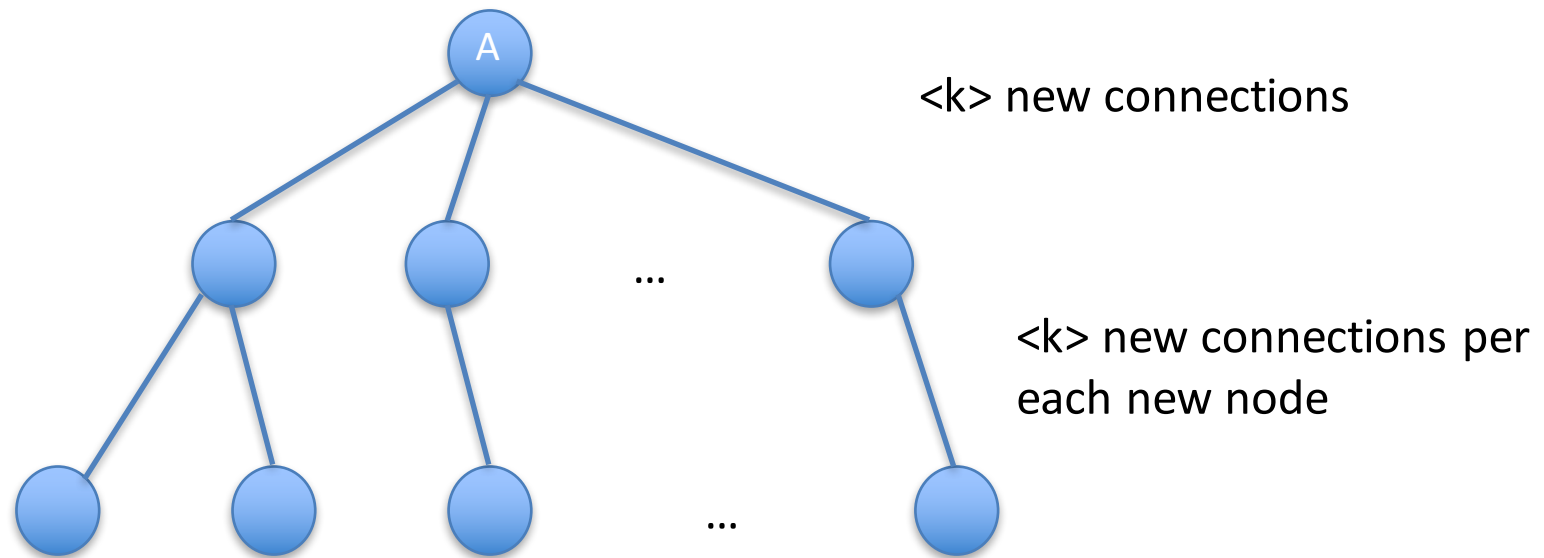


- The **diameter** of a random graph and the **average path length** of the graph have been demonstrated to be:

$$d = \frac{\ln(N)}{\ln(pN)} = \frac{\ln(N)}{\ln(\langle k \rangle)} \approx l_{rand}$$

The average distance between two nodes is quite small wrt to the size of the graph.

# Random Graph Diameter: An Intuition



- The nodes at distance  $l$  from  $A$  will be  $\sim \langle k \rangle^l$

$$N = k^l$$
$$\log N = l \log k$$
$$l = \log N / \log k$$



# Relationship of $\langle k \rangle$ and connectivity

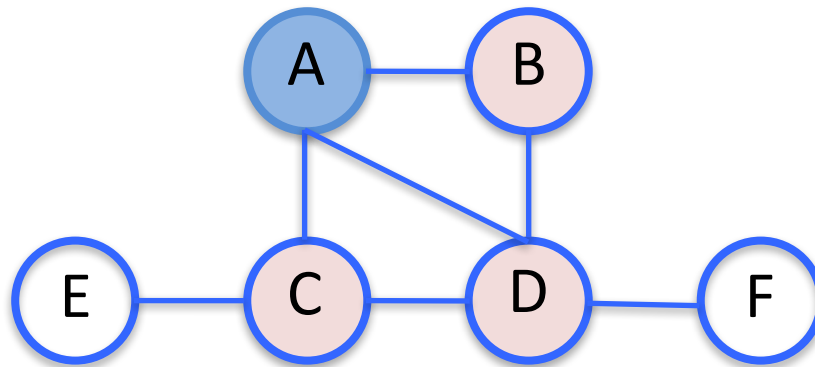
---

- $\langle k \rangle = \text{average degree (np)}$
- If  $\langle k \rangle < 1$  disconnected network
- If  $\langle k \rangle > 1$  a giant component appears
- If  $\langle k \rangle \geq \ln(N)$  graph is totally connected

# Clustering Coefficient



- The **clustering coefficient** defines the proportion of  $A$ 's neighbours ( $N(A)$ ) which are connected by an edge (are friends).
- The number of triangles in which  $A$  is involved wrt to the ones it could be involved in.



# Formally: Clustering Coefficient

---

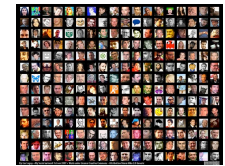


Local Clustering Coefficient

$$C_i = \frac{2|\{e_{jk}\}|}{k_i(k_i - 1)} : v_j, v_k \in N_i, e_{j,k} \in E$$

Network Clustering Coefficient

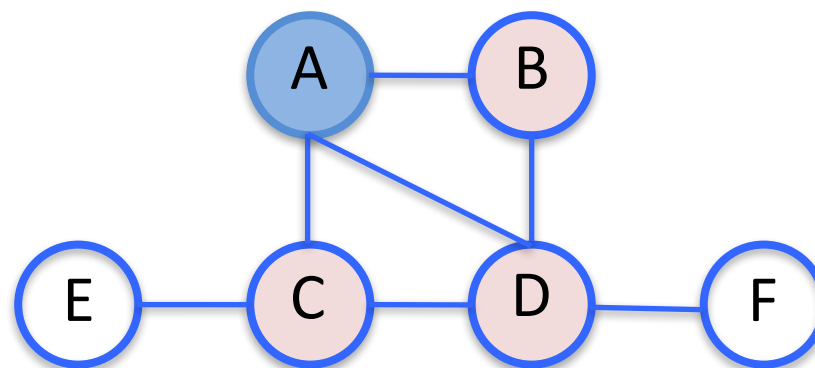
$$CG = \frac{1}{N} \sum_i C_i$$



# Clustering Coefficient: Example

---

- $K_i=3$
- Nominator=  $2*2=4$
- Denominator=6
- $C_i=4/6=2/3$



# Clustering Coefficient of a Random Graph

---



- The clustering coefficient of a random graph is

$$C_{rand} = p = \frac{\langle k \rangle}{N}$$

- The probability that 2 neighbours of a node are connected is equal to the probability that 2 random nodes are connected:
  - $p * k_i(k_i-1)$  where  $k_i(k_i-1)$  is the total number of pairs
  - Then the Clustering is  $p * k_i(k_i-1)/k_i(k_i-1)$
- Is this mirroring the clustering coefficient of real networks?

# Question

---



- Are Random Graphs representatives of Real Networks?





# Summary

---

- We have introduced graphs definitions and measures.
- Random graphs are a first examples of models for networks.

# References



- Material from Chapter 1, 2 of
  - **D. Easley, J. Kleinberg. Networks, Crowds, and Markets: Reasoning About a Highly Connected World. Cambridge University Press, 2010.**
- P. Sheridan Dodds, R. Muhamad, and D. J. Watts. An Experimental Study of Search in Global Social Networks. *Science* 8. August 2003: 301 (5634), 827-829.
- J. Ugander, B. Karrer, L. Backstrom, C. Marlow. The Anatomy of the Facebook Social Graph, <http://arxiv.org/abs/1111.4503>
- L. Backstrom, P. Boldi, M. Rosa, J. Ugander, S. Vigna. Four Degrees of Separation, <http://arxiv.org/abs/1111.4570>
- R. Albert, A. Barabasi. Statistical Mechanics of Complex Networks. *Reviews of Modern Physics* (74). Jan. 2002.
- S. Boccaletti, V. Latora, Y. Moreno, M. Chavez, D.-U. Hwang, *Complex Networks: Structure and Dynamics Physics Reports* 424 (2006) 175 .