

Random Networks

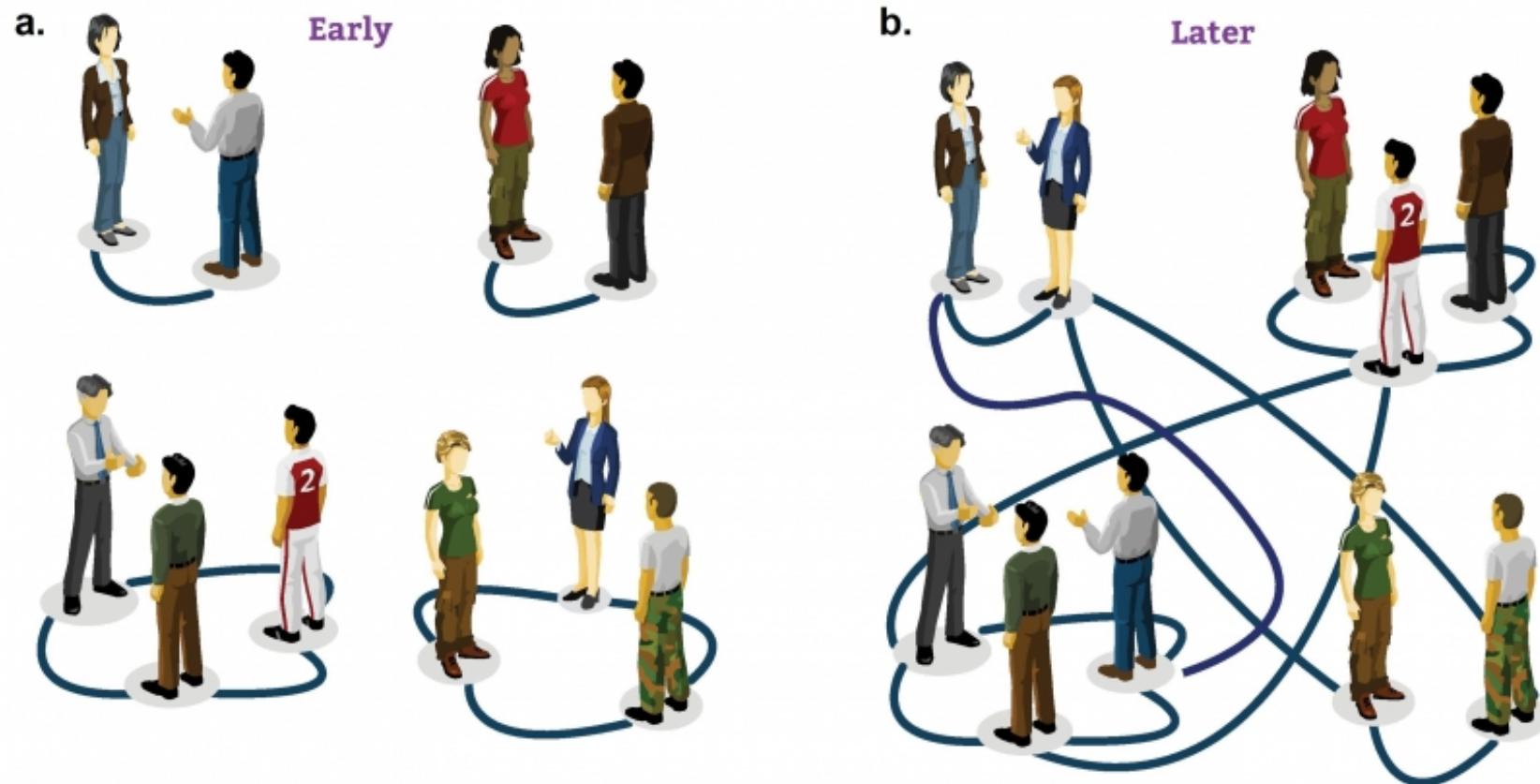
Prof. Pascal Frossard
Signal Processing Laboratory (LTS4)

Some slides are inspired from Prof. Barabási's class on Network Science (www.BarabasiLab.com)

Outline

- *Why network models?*
- Random network model
- Random network evolution
- Small worlds
- Clustering coefficient

Will we run out of good wine?



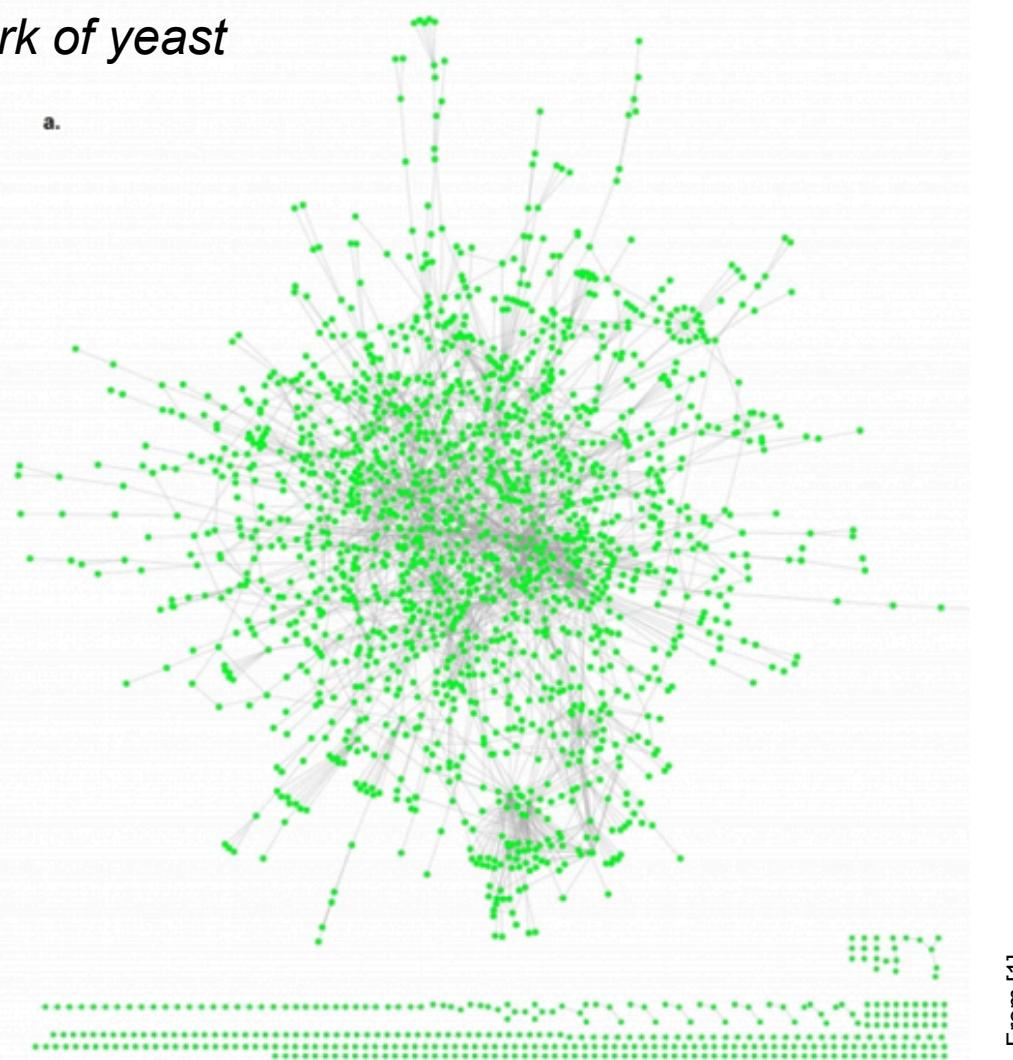
From [1]

Why network models?

- Most (simple and tractable) network models do not describe accurately the real data :(
- However, network models are very useful to understand real networks and their properties
 - It will help us calculate many quantities, that can then be compared to the real data, understanding to what degree a particular property is valid
- In order to identify these properties, we need to understand how a network would look like if it is driven entirely by a model
- In particular, the random network model is a sort of reference or benchmark model

Is that a random network?

protein interaction network of yeast

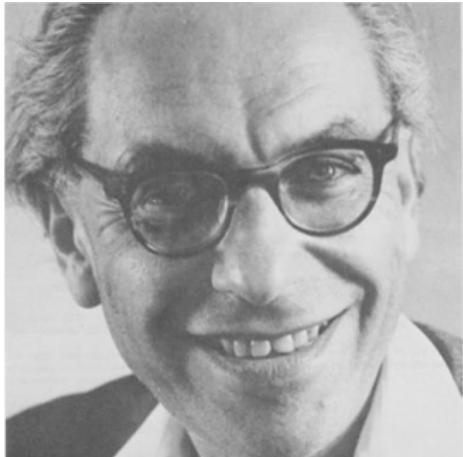


Outline

- Why network models?
- *Random network model*
- Random network evolution
- Small worlds
- Clustering coefficient

Erdös-Rényi model (1960)

Pál Erdös
(1913-1996)



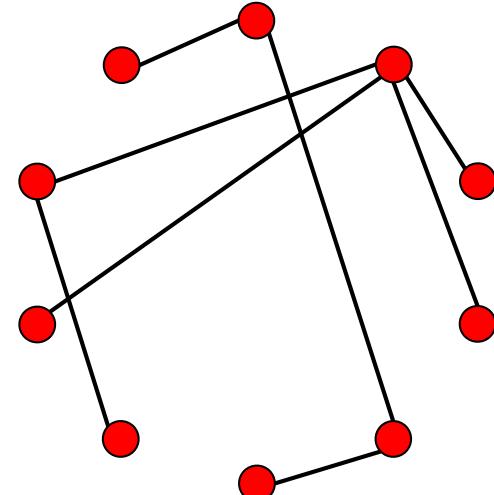
Alfréd Rényi
(1921-1970)

$G(N, L)$ model: N labeled nodes are connected with L randomly placed links.

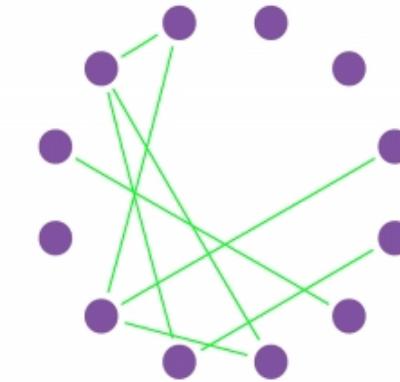
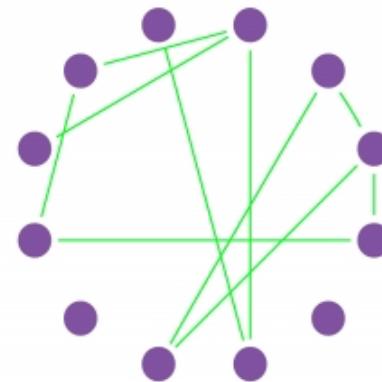
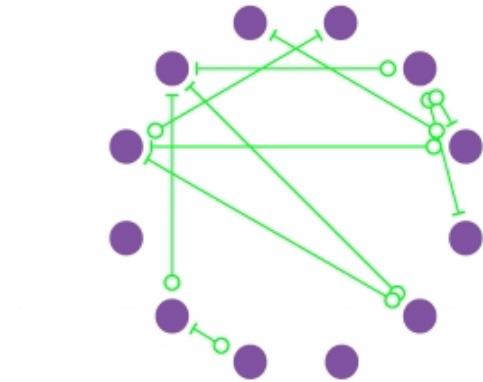
- Gilbert's $G(N, p)$ model: a random network model is a network where each pair of nodes is connected with probability p
 - Example: $p = 1/6$

$$N = 10$$

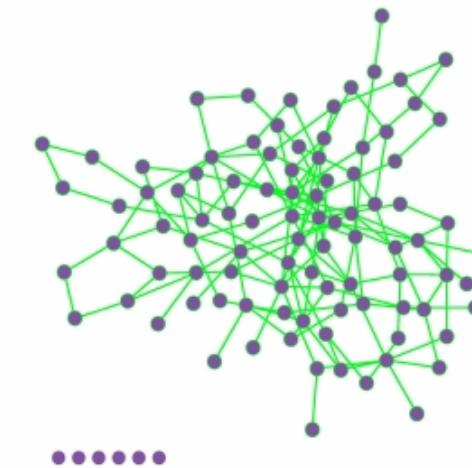
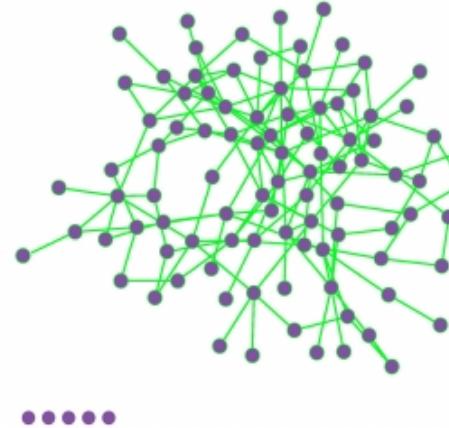
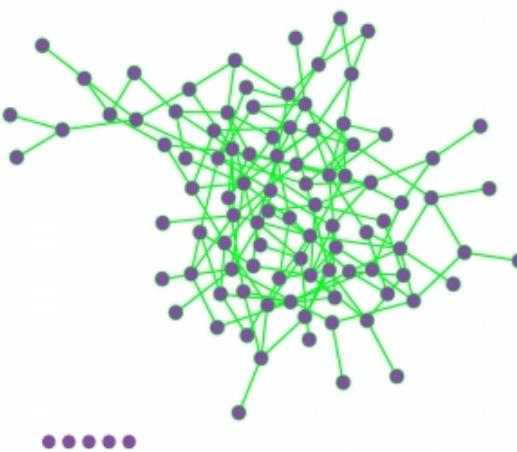
$$\langle k \rangle \sim 1.5$$



Random network examples



$p = 1/6$
 $N = 12$



$p = 0.03$
 $N = 100$

From [1]

Number of links in a $G(N,p)$ network

- Probability for the N -node network to have L links (binomial distribution):

$$p_L = \binom{\frac{N(N-1)}{2}}{L} p^L (1-p)^{\frac{N(N-1)}{2}-L}$$

Probability to have L links

Number of ways to place L linksProbability that other attempts did not result in a link

- Expected number of links in the random graph

$$\langle L \rangle = \sum_{L=0}^{\frac{N(N-1)}{2}} L p_L = p \frac{N(N-1)}{2}$$

Maximum number of links

- Average degree

$$\langle k \rangle = \frac{2 \langle L \rangle}{N} = p(N-1)$$

Maximum number of links for one node

Degree distribution in $G(N,p)$

- Degree distribution for a random network (binomial distribution)

$$p_k = \binom{N-1}{k} p^k (1-p)^{N-1-k}$$

Number of ways to select k links

- Degree distribution for $\langle k \rangle \ll N$ (sparse networks) is approximated by the one of the Poisson distribution

$$p_k = e^{-\langle k \rangle} \frac{\langle k \rangle^k}{k!}$$

Simpler, but does not depend on N and valid only for sparse networks!

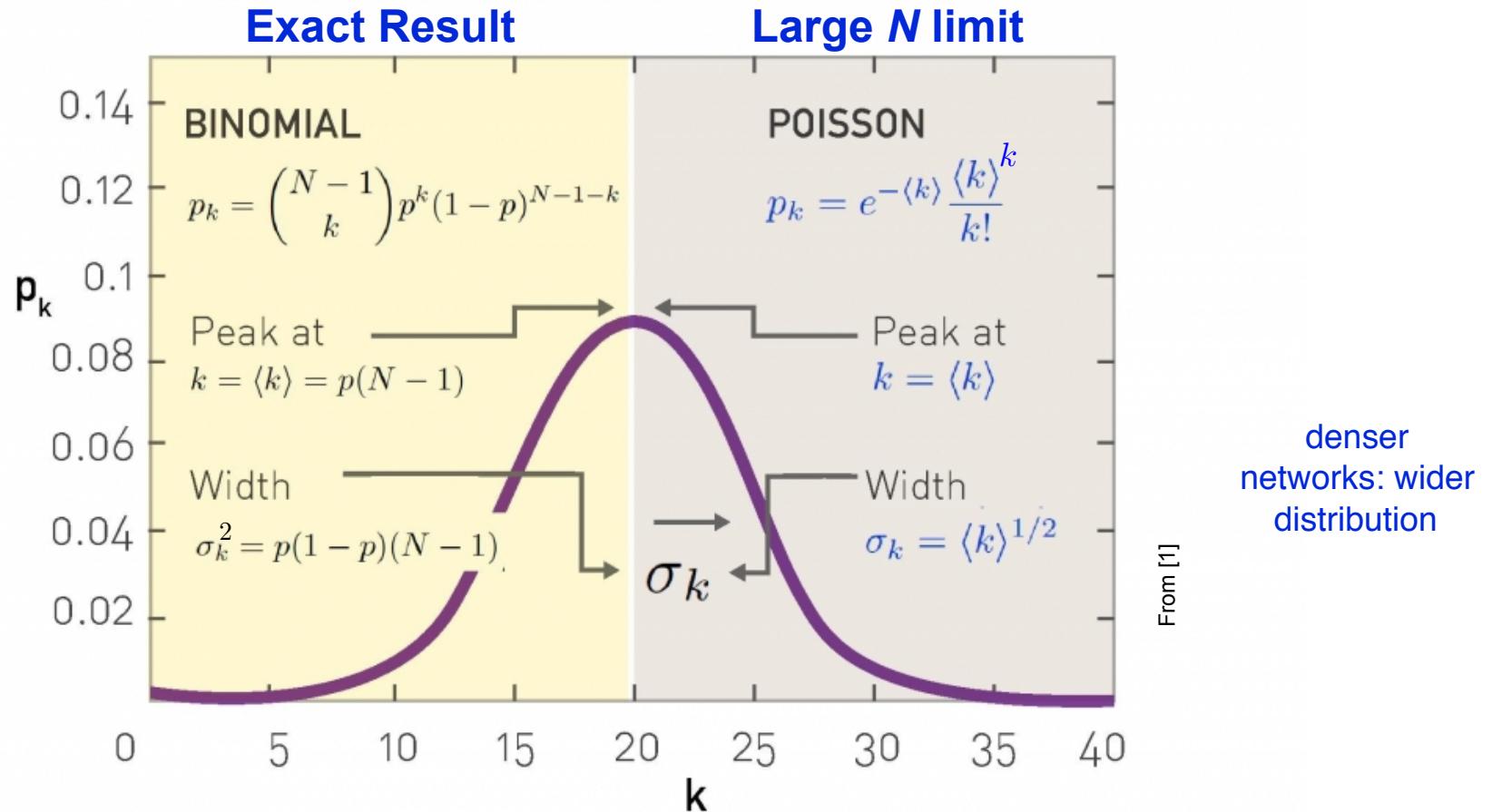
From [1]

Poisson distribution

- Exact form of the distribution is binomial: $p_k = \binom{N-1}{k} p^k (1-p)^{N-1-k}$
- But $\binom{N-1}{k} = \frac{(N-1)(N-1-1)(N-1-2)...(N-1-k+1)}{k!} \approx \frac{(N-1)^k}{k!}$ if $k \ll N$
- And $\ln[(1-p)^{(N-1)-k}] = (N-1-k) \ln\left(1 - \frac{\langle k \rangle}{N-1}\right)$ $\textcolor{blue}{p = \frac{\langle k \rangle}{N-1}}$
with $\ln(1+x) = \sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n} x^n = x - \frac{x^2}{2} + \frac{x^3}{3} - \dots, \forall |x| \leq 1$ (series expansion)
- So that $\ln[(1-p)^{(N-1)-k}] \approx (N-1-k) \frac{-\langle k \rangle}{N-1} = -\langle k \rangle \left(1 - \frac{k}{N-1}\right) \approx -\langle k \rangle$
and $(1-p)^{N-1-k} = e^{-\langle k \rangle}$ valid for $k \ll N$
- Finally $p_k = \binom{N-1}{k} p^k (1-p)^{(N-1)-k} = \frac{(N-1)^k}{k!} p^k e^{-\langle k \rangle} = \frac{(N-1)^k}{k!} \left(\frac{\langle k \rangle}{N-1}\right)^k e^{-\langle k \rangle}$
or
$$p_k = e^{-\langle k \rangle} \frac{\langle k \rangle^k}{k!}$$
 (Poisson distribution approximation)

Degree distribution in $G(N,p)$

p increases:
peak moves to
the right,
network gets
denser

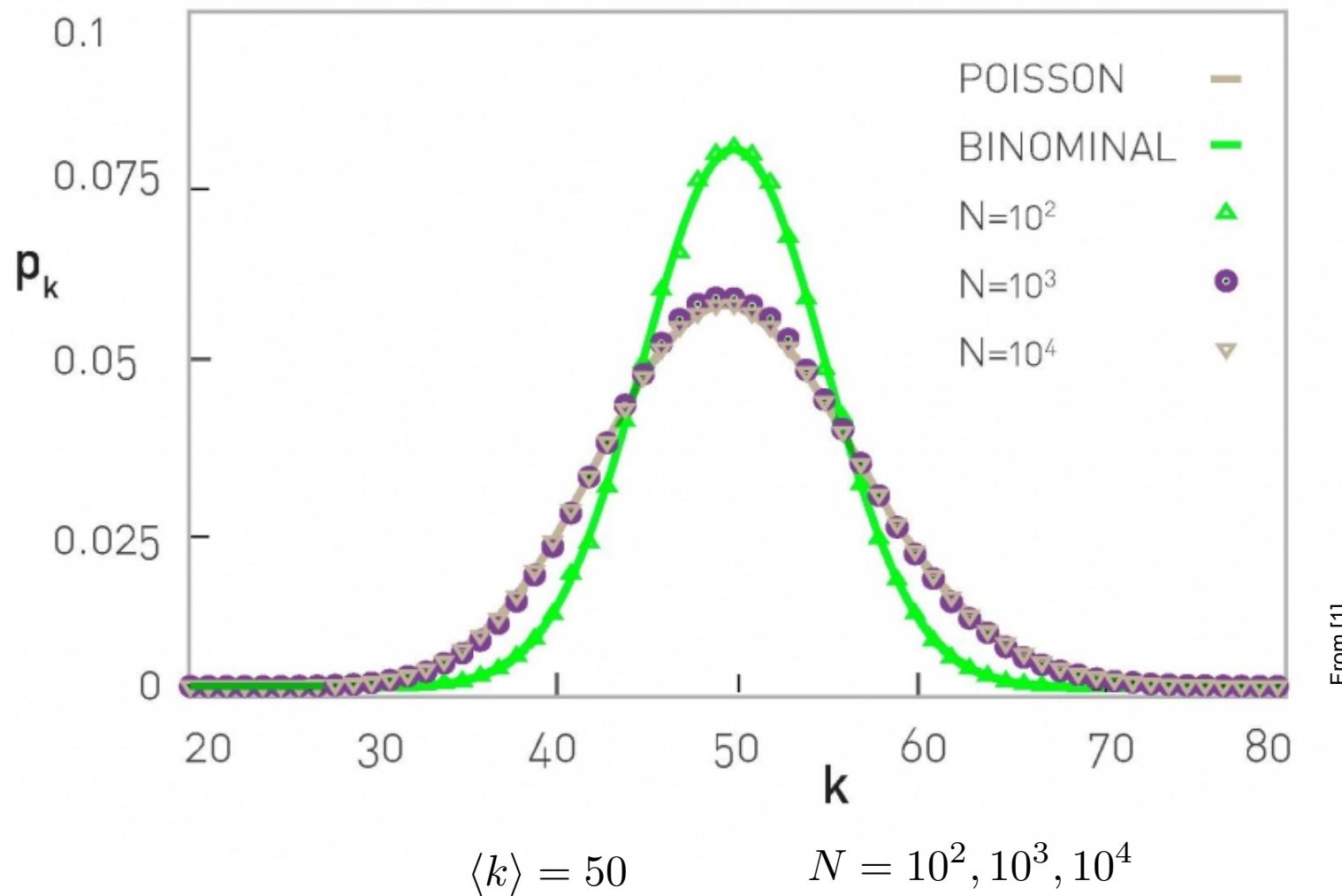


We further have

$$\frac{\sigma_k}{\langle k \rangle} = \left[\frac{1-p}{p} \frac{1}{(N-1)} \right]^{\frac{1}{2}} \approx \frac{1}{(N-1)^{\frac{1}{2}}}$$

N increases: distribution gets more narrow

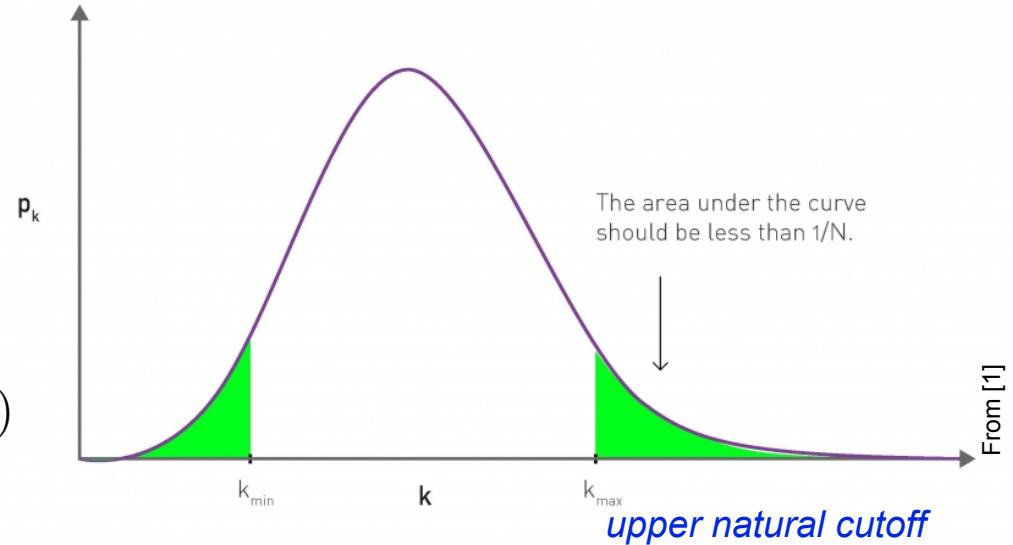
Network size influence



Minimum and maximum degrees

- Maximum (resp. minimum) degree is chosen so that there remains at most one node with a higher (resp. lower) degree value

- the area under the curve is $1 - P(k_{\max})$
cumulative degree distribution



- then we should have $N[1 - P(k_{\max})] \approx 1$
- For the Poisson distribution:

$$1 - P(k_{\max}) = 1 - e^{-\langle k \rangle} \sum_{k=0}^{k_{\max}} \frac{\langle k \rangle^k}{k!} \approx e^{-\langle k \rangle} \sum_{k=k_{\max}+1}^{\infty} \frac{\langle k \rangle^k}{k!} \approx e^{-\langle k \rangle} \frac{\langle k \rangle^{k_{\max}+1}}{(k_{\max} + 1)!}$$

$\approx 1/N$

- Similarly, for the minimum degree, we should have $NP(k_{\min} - 1) \simeq 1$
- For the Poisson distribution: $P(k_{\min} - 1) = e^{-\langle k \rangle} \sum_{k=0}^{k_{\min}-1} \frac{\langle k \rangle^k}{k!}$

approximation with largest term in the sum

Poisson vs social network

- Assume that a typical person knows about 1000 people, we have

$$\langle k \rangle = 1000 \quad \text{and} \quad N = 7 \times 10^9$$

- With a Poisson distribution (random network), we would have

$$k_{max} = 1,185 \quad \text{and} \quad k_{min} = 816$$

that are very similar values. The dispersion of the degree value is also very narrow,

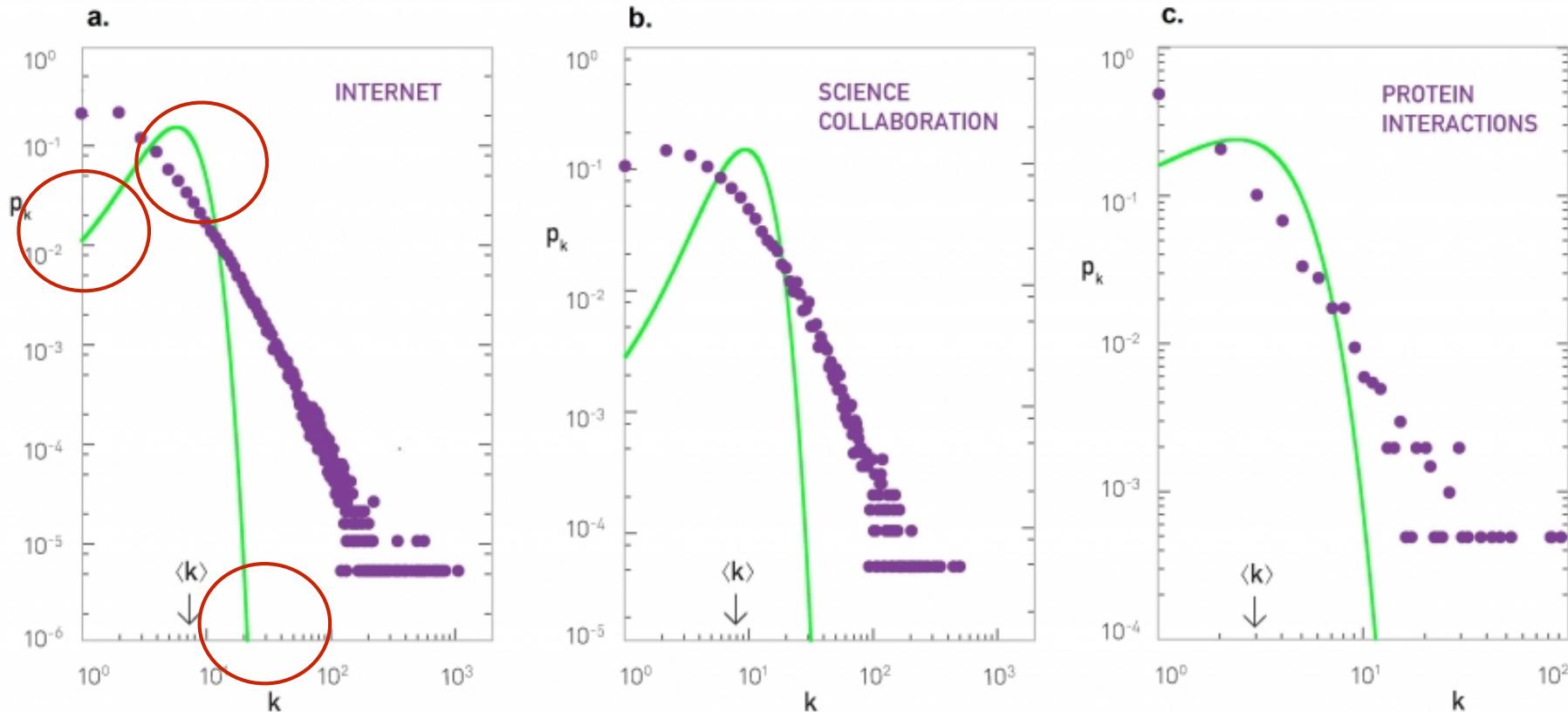
$$\sigma_k = \langle k \rangle^{\frac{1}{2}} = 31.62$$

- In a large random network, the degree of most nodes is close to $\langle k \rangle$. This is clearly not representative of social networks, with nodes with much higher degrees!
- Hubs are missing in random networks!

$$p_k = e^{-\langle k \rangle} \frac{\langle k \rangle^k}{k!} \quad \xrightarrow{\text{Stirling approximation}} \quad p_k = \frac{e^{-\langle k \rangle}}{\sqrt{2\pi k}} \left(\frac{e \langle k \rangle}{k} \right)^k$$

Both factors decay fast for large k !

Real networks are not Poisson



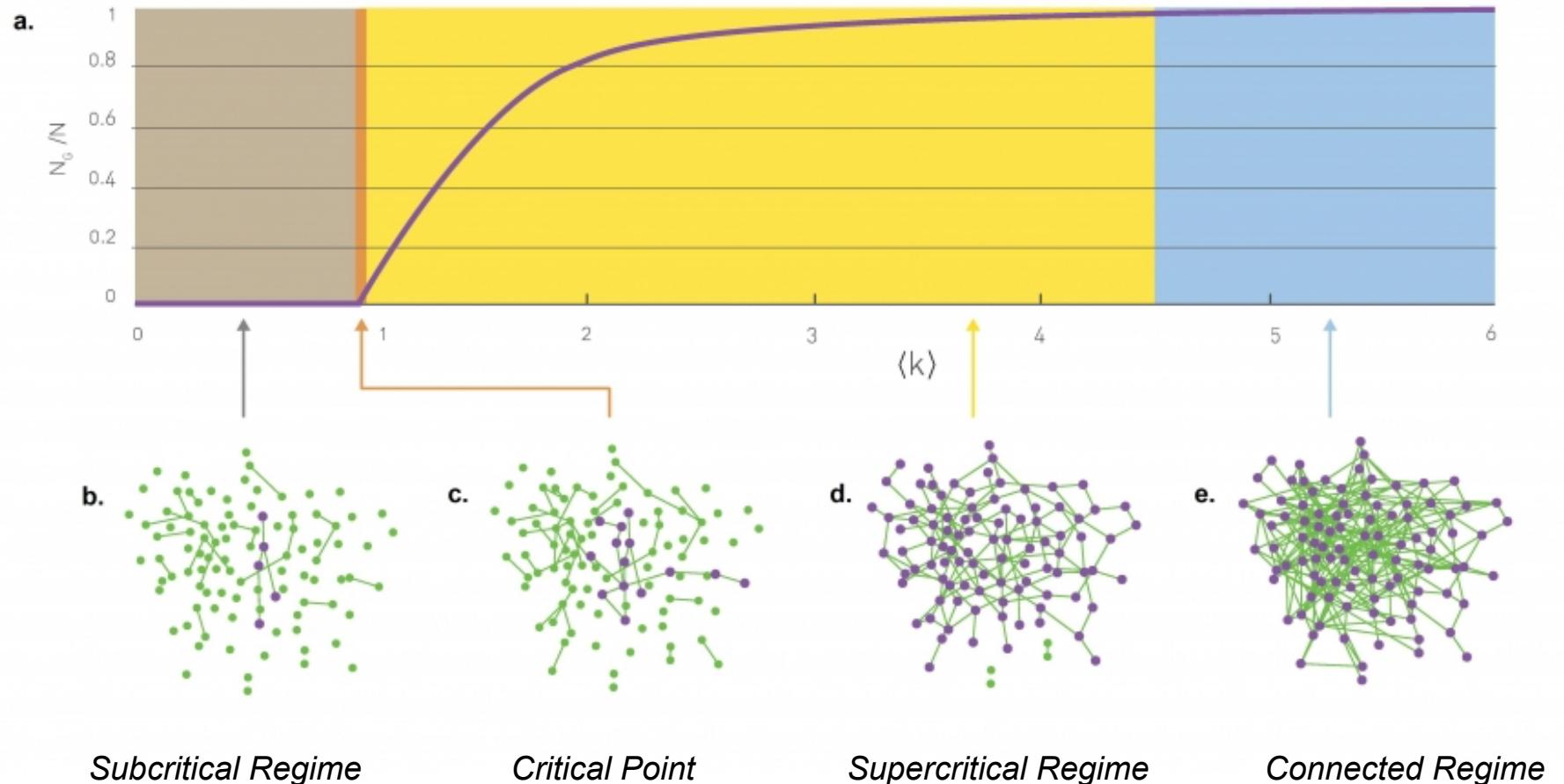
From [1]

Poisson prediction (green curves) are not good fits!

Outline

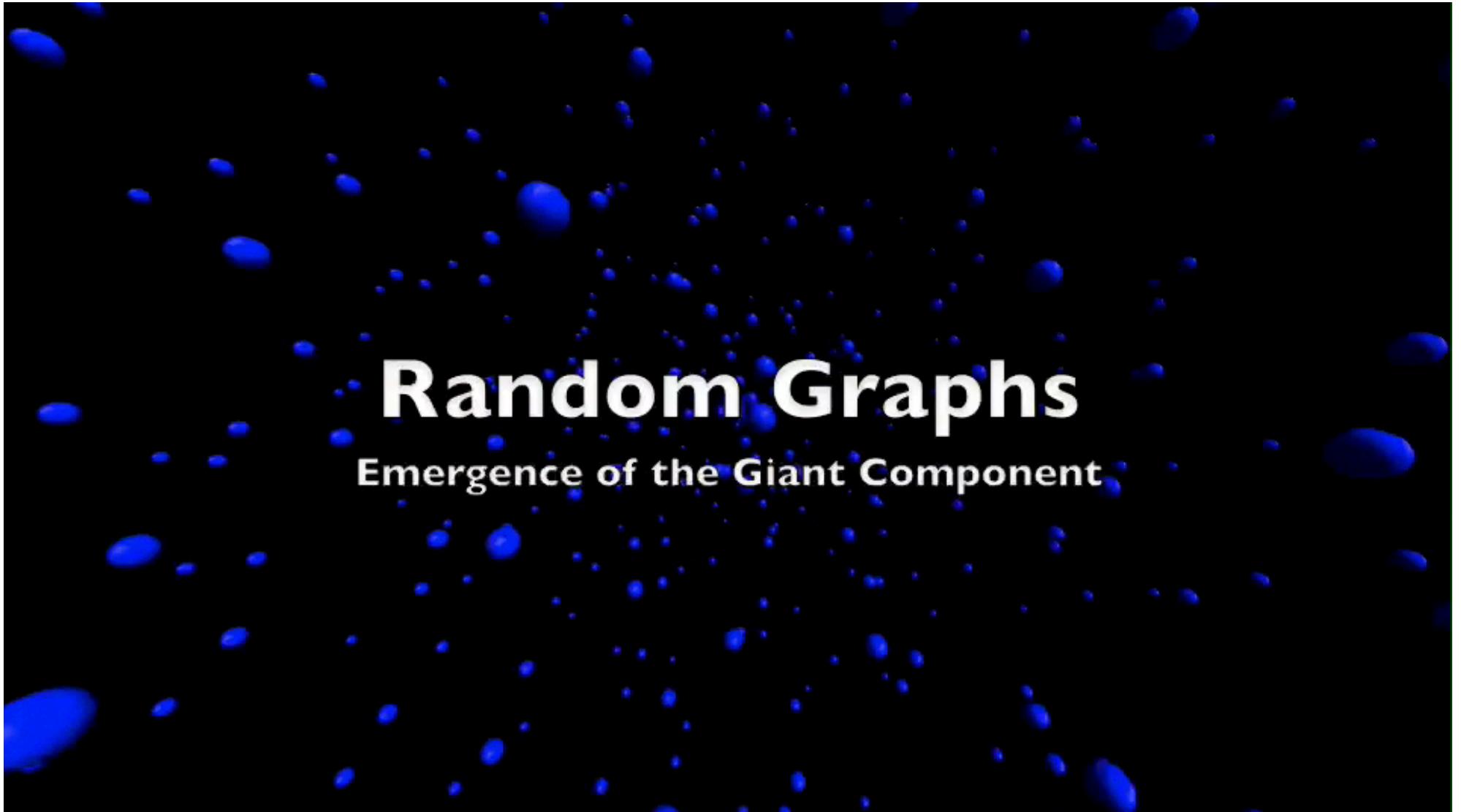
- Why network models?
- Random network model
- *Random network evolution*
- Small worlds
- Clustering coefficient

Evolution of random networks



How does the network change?

Emergence of giant components



Emergence of a giant component

- We look at the size of the giant component, N_G
- First, there are two extreme cases:
 - If $p = 0$, we have $\langle k \rangle = 0$, $N_G = 1$ and $\frac{N_G}{N} \rightarrow 0$ for large N
 - If $p = 1$, we have $\langle k \rangle = N - 1$, $N_G = N$ and $\frac{N_G}{N} = 1$
- What happens in between might look surprising
 - it is not a gradual increase of the size of the giant component with increasing p
 - there is rather the rapid emergence of a large cluster after some *critical value*
 - the ‘phase transition’ happens *already* when each node has more than one link, i.e., after

$$\langle k_c \rangle = 1$$

[Erdős and Rényi, 1959]

- equivalently, as $\langle k \rangle = p(N - 1)$, we also have

$$p_c = \frac{1}{N - 1} \approx \frac{1}{N}$$

The larger the network, the smaller p

Giant component (GC)

- Fraction of nodes not in GC (of size N_G): $u = 1 - \frac{N_G}{N}$
 - for node i not to be part of GC, it has either no link with node j in the GC (probability $1-p$) or it has a link with j , which is however not in the GC (probability pu)
 - hence, the probability that i is not linked to the GC via any other node j is

$$u = (1 - p + pu)^{N-1}$$

- Combining both relations, and using $p = \frac{\langle k \rangle}{N-1}$, we have

$$\ln u = (N-1) \ln \left[1 - \frac{\langle k \rangle}{N-1} (1-u) \right] \approx (N-1) \left[-\frac{\langle k \rangle}{N-1} (1-u) \right] = -\langle k \rangle (1-u)$$

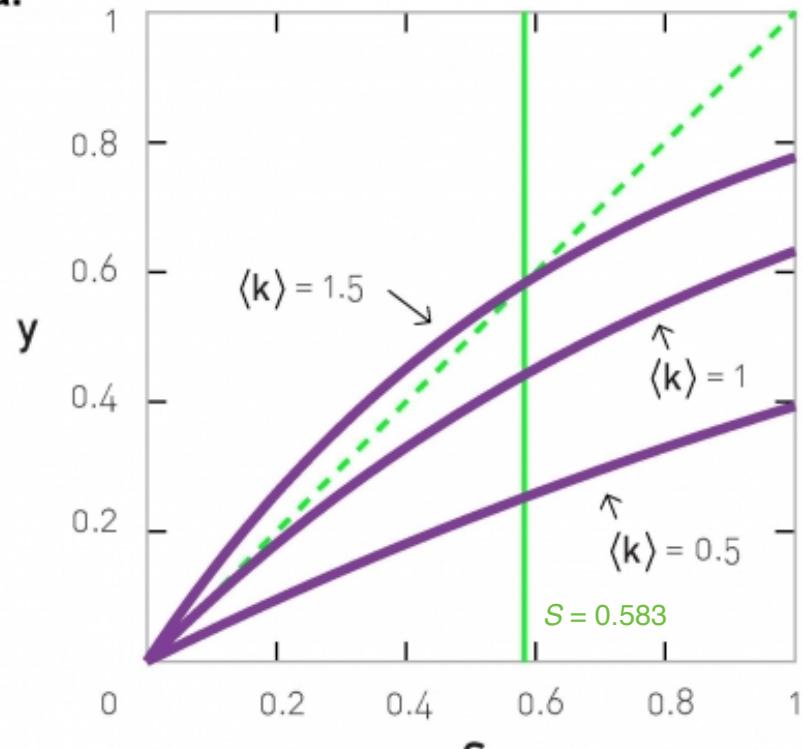
series expansion for $\ln(1+x)$ $\langle k \rangle \ll N$

- Equivalently, $u = \exp[-\langle k \rangle (1-u)]$
- Finally, if we rather look at the fraction of nodes in GC, $S = \frac{N_G}{N}$, we have $S = 1 - u$, or

$$S = 1 - e^{-\langle k \rangle S}$$

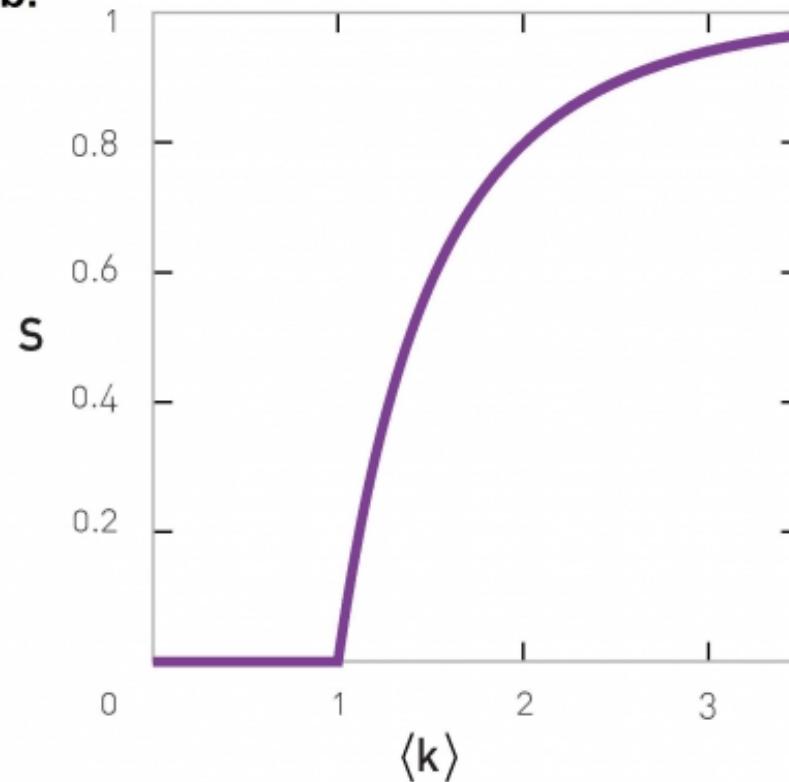
Evolution of GC size

a.



$$y = 1 - e^{-\langle k \rangle S}$$

b.



From [1]

$$\text{Solution of } S = 1 - e^{-\langle k \rangle S}$$

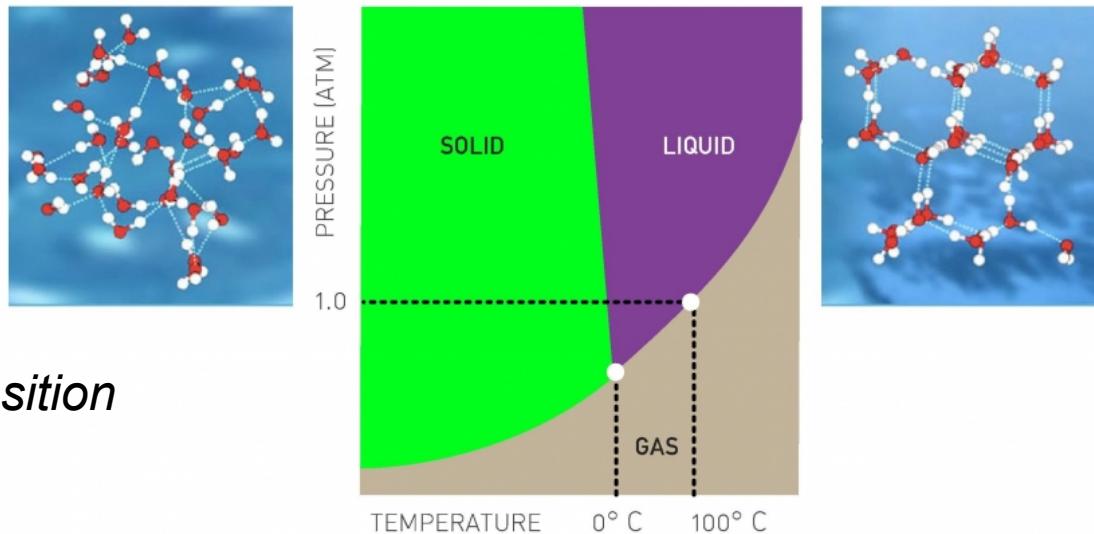
First non-trivial solution happens when, by changing $\langle k \rangle$, there is intersection with $y=S$ at $S=0$. Namely:

$$\frac{d}{dS} (1 - e^{-\langle k \rangle S}) = 1$$

$$\langle k \rangle e^{-\langle k \rangle S} = 1 \quad \text{or...} \quad \langle k \rangle = 1 \quad [\text{Erdős and Rényi, 1959}]$$

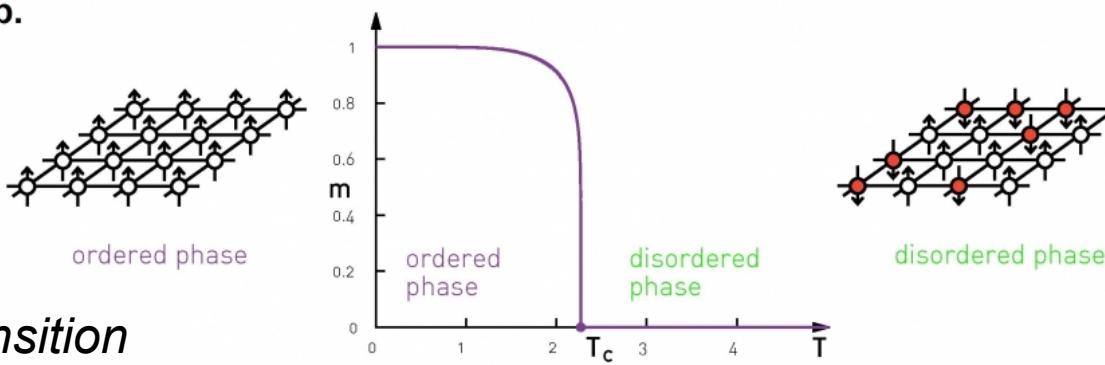
Other phase transition examples

a.



water-ice phase transition

b.



magnetic phase transition

From [1]

Component size distribution

- Probability that a random node belongs to a component of size s (not GC)

$$p_s \sim \frac{(s \langle k \rangle)^{s-1}}{s!} e^{-\langle k \rangle s}$$

- As $\langle k \rangle^{s-1}$ can be replaced by $\exp[(s-1)\ln\langle k \rangle]$ and since $s! = \sqrt{2\pi s} \left(\frac{s}{e}\right)^s$

Stirling approximation, for large s

$$p_s \sim s^{-3/2} e^{-(\langle k \rangle - 1)s + (s-1)\ln\langle k \rangle}$$

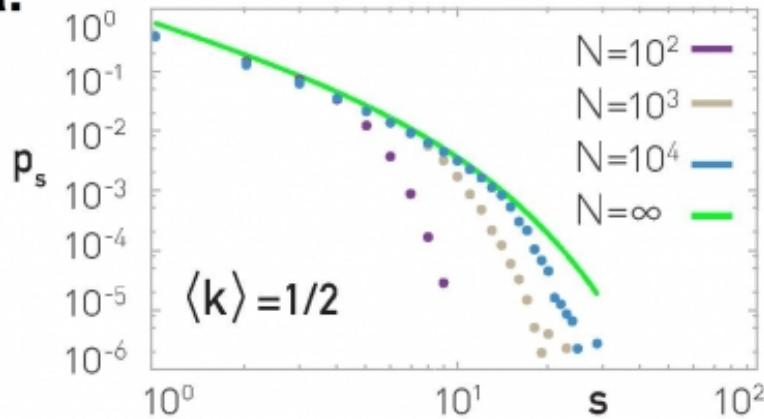
Dominated by the quickly decreasing exponential, for large s

- large components (diff than GC) are prohibited as p_s goes quickly to 0 for large s
- At the critical point $\langle k \rangle = 1$, we however have $p_s \sim s^{-3/2}$
 - power-law, with clusters of widely different sizes at the critical point !

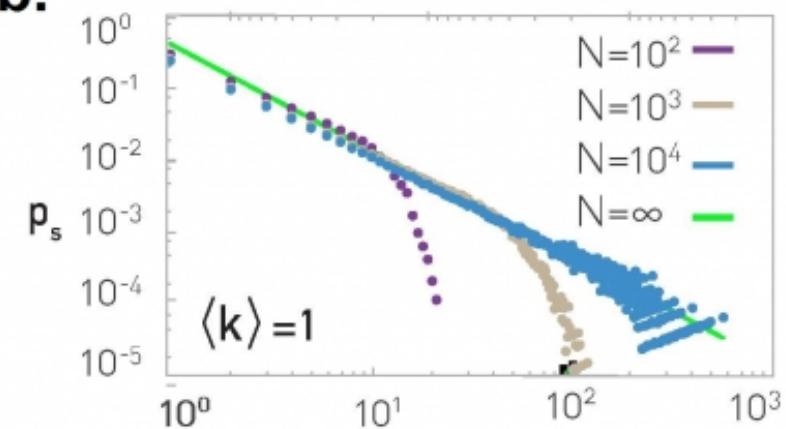
See [2] for details

Component size illustration

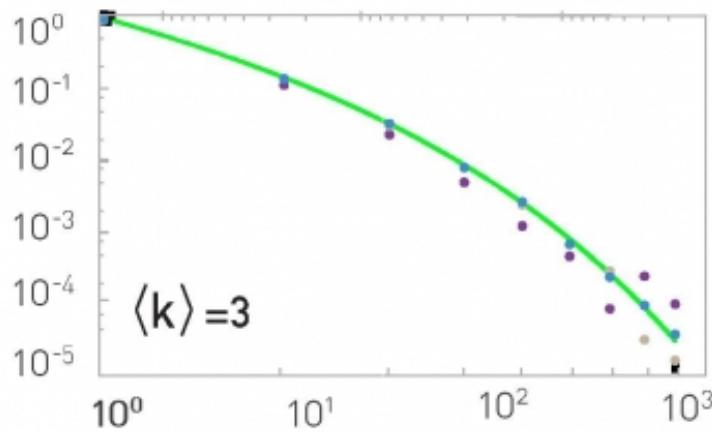
a.



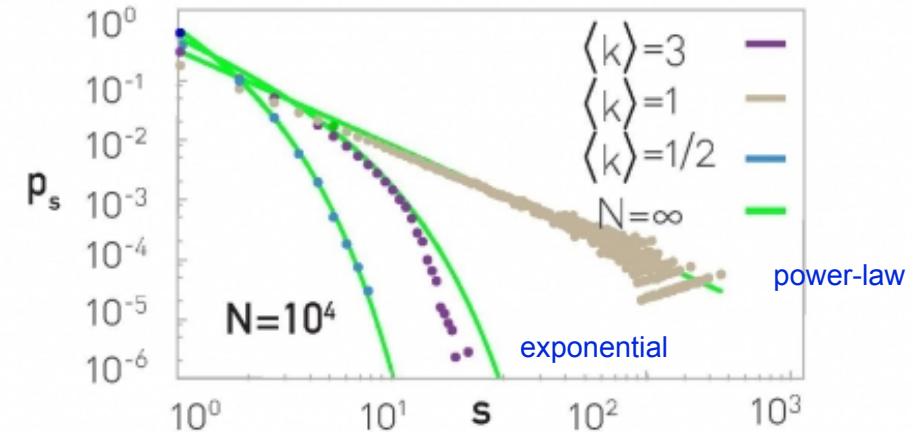
b.



c.



d.



From [1]

Average component size

- The average component size (excluding GC) is given as

$$\langle s \rangle = \frac{1}{1 - \langle k \rangle + \langle k \rangle N_G/N}$$

- for $\langle k \rangle < 1$, $N_G = 0$ and $\langle s \rangle = \frac{1}{1 - \langle k \rangle}$

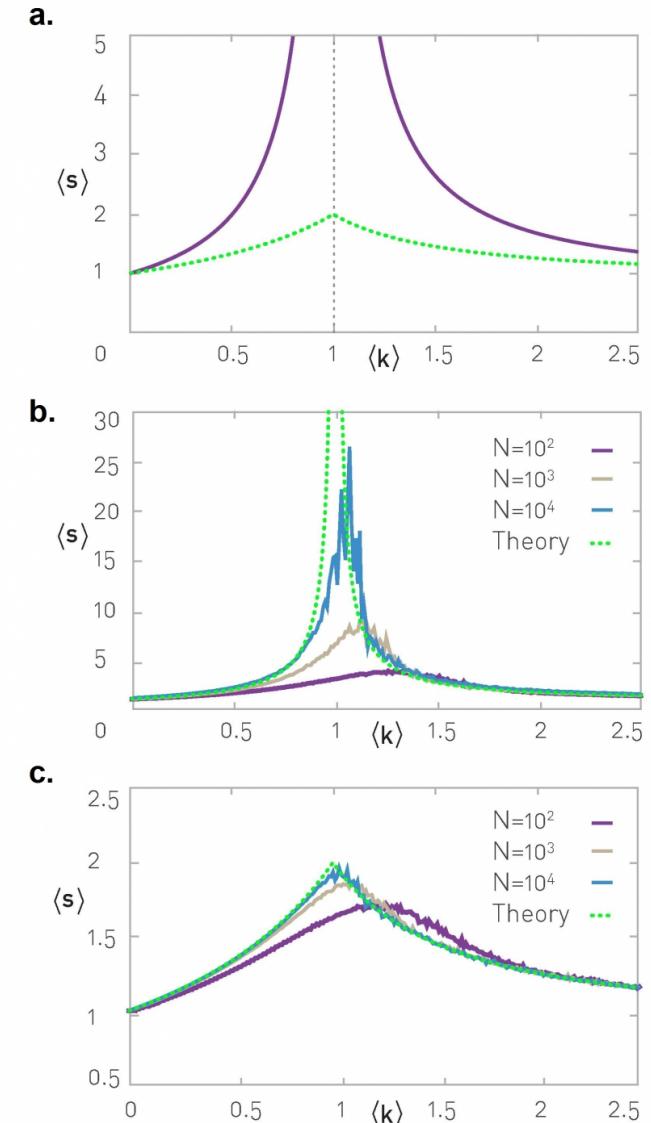
Divergence, when $\langle k \rangle = 1$

- There is a bias in picking a random node - chances to belong to a larger cluster is higher. Correcting the bias (linear in s), we have

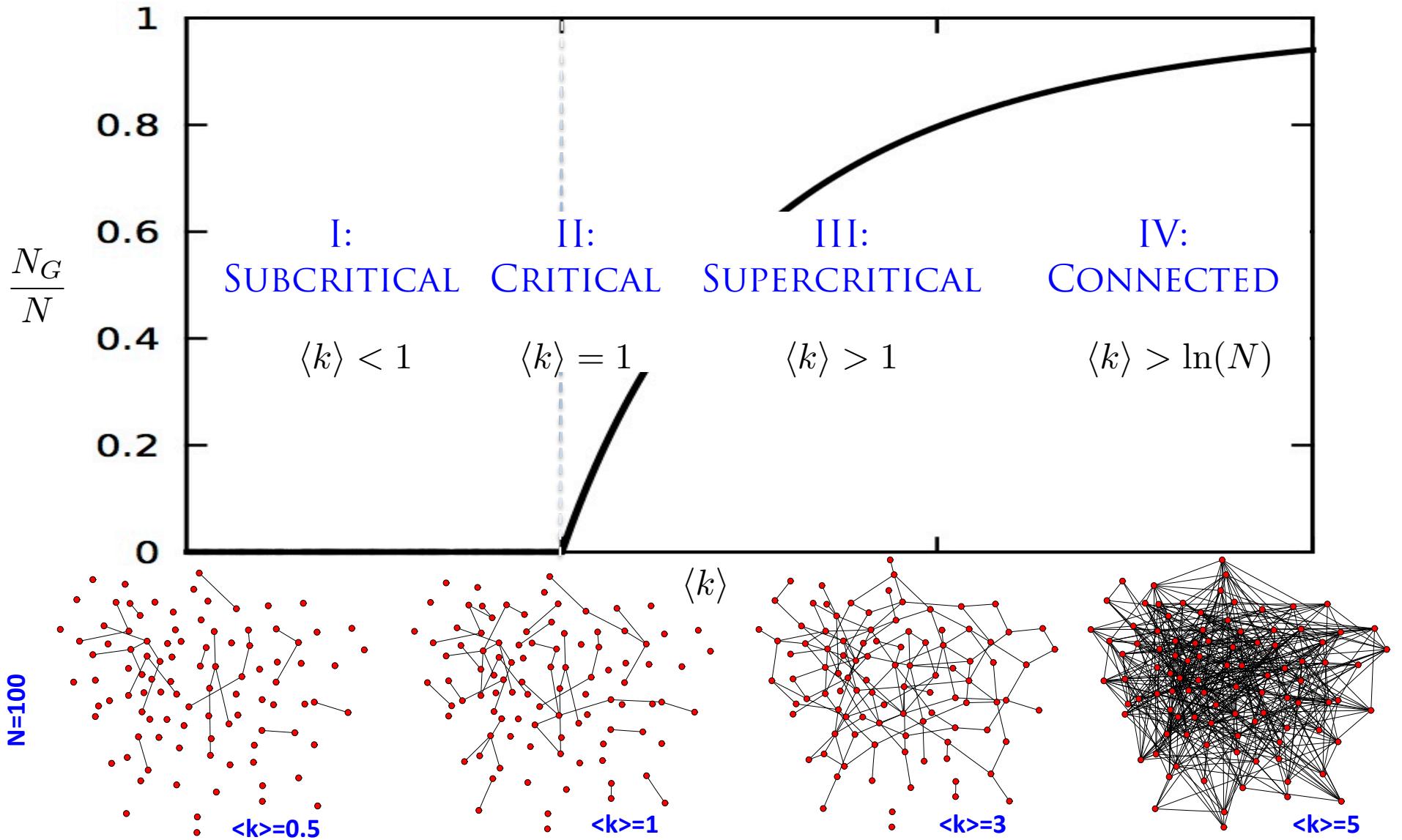
$$\langle s' \rangle = \frac{2}{2 - \langle k \rangle + \langle k \rangle N_G/N}$$

It corresponds to inspecting each cluster one by one

See [2] for details



Regimes



Subcritical regime

- The subcritical regime corresponds to the lowest average degree, i.e., $0 \leq \langle k \rangle < 1$ or $p < 1/N$
- For $\langle k \rangle = 0$, the network consists of N isolated nodes
- Increasing $\langle k \rangle$ corresponds to adding $N\langle k \rangle/2 = pN(N - 1)/2$ links
 - the number of links stays small in the regime, and there are mainly tiny clusters
- The largest cluster is a tree with size $N_G \approx \ln(N)$ and for large N
$$N_G/N = \ln N/N \rightarrow 0$$

Numerous tiny components of comparable size - no giant component!

Critical regime

- Critical point $\langle k \rangle = 1$ (or $p = 1/N$) between regimes where there is no GC, and where there is a GC, respectively.

- The relative size of the largest component is still 0:

$$N_G \sim N^{2/3} \text{ so that } N_G/N \sim N^{-1/3} \rightarrow 0 \text{ for large } N$$

- Still, there is a significant jump in GC size, e.g.,

$N = 1000$	$\ln(N) \sim 6.9$	$N^{2/3} \sim 95$
$N = 7 \times 10^9$	$\ln(N) \sim 22$	$N^{2/3} \sim 3,659,250$

- Yet, at the critical point the largest component contains only a vanishing fraction of the total number of nodes in the network.

Numerous small components of rather different sizes coexist.

Supercritical regime

- This regime has the most relevance to real systems - for the first time it has a giant component that looks like a network.
- Close to the critical point $\langle k \rangle = 1$ we have

$$\frac{N_G}{N} \sim \langle k \rangle - 1 \quad N_G \sim (p - p_c)N$$

- The giant component contains a finite fraction of the nodes that grows with $\langle k \rangle$
- For large $\langle k \rangle$ the size of the GC grows non-linearly with $\langle k \rangle$

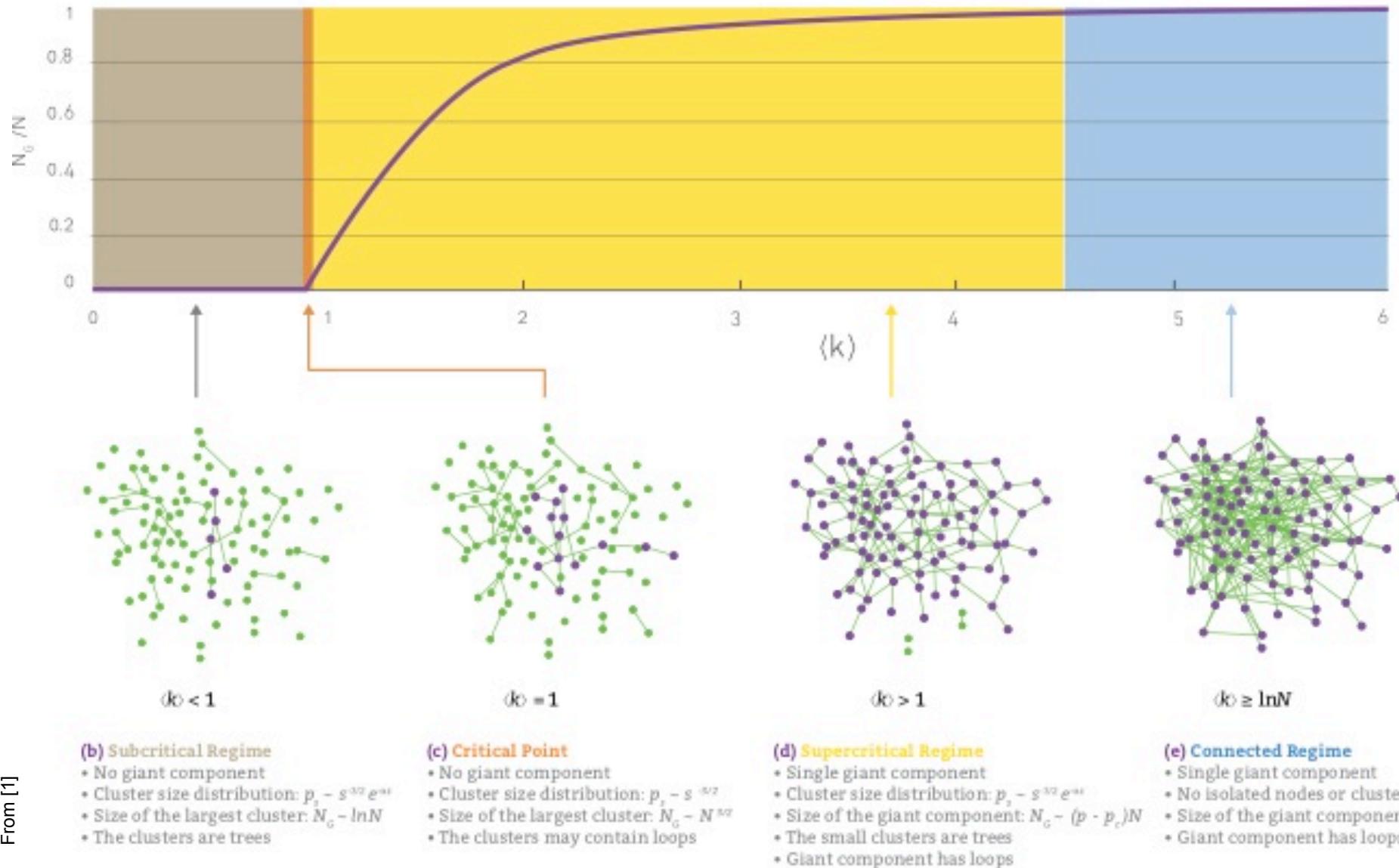
Numerous isolated components with a giant component.

Connected regime

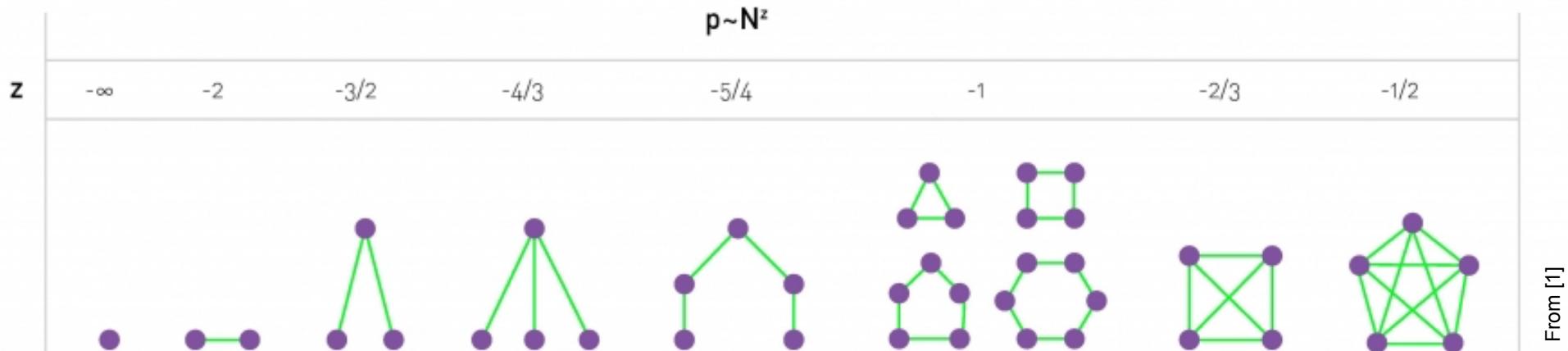
- For large p , all nodes are absorbed by the GC and $N_G = N$
 - the network becomes connected
 - the network starts relatively sparse, and turns into a complete graph at $\langle k \rangle = N - 1$
- The threshold after which most nodes are part of GC is $\langle k \rangle > \ln N$ or $p > \ln(N)/N$
 - probability that a random node is not connected to GC: $(1 - p)^{N_G} \approx (1 - p)^N$
 - expected number of isolated nodes: $I_N = N(1 - p)^N = N \left(1 - \frac{N \cdot p}{N}\right)^N \approx Ne^{-Np}$
$$(1 - x/n)^n \approx e^{-x}$$
 - for $I_N = 1$, p needs to satisfy $Ne^{-Np} = 1$, hence $p = \frac{\ln N}{N}$

One single giant component.

Recap: random network evolution



Network evolution in graph theory



- In the random graph theory literature, the connection probability $p(N)$ scales with N^z with $z \in]-\infty, 0]$
- The threshold probabilities at which different subgraphs appear is governed by z .

R. Albert and A.-L. Barabási. Statistical mechanics of complex networks. *Reviews of Modern Physics*, 74:47-97, 2002.

Real Networks are supercritical

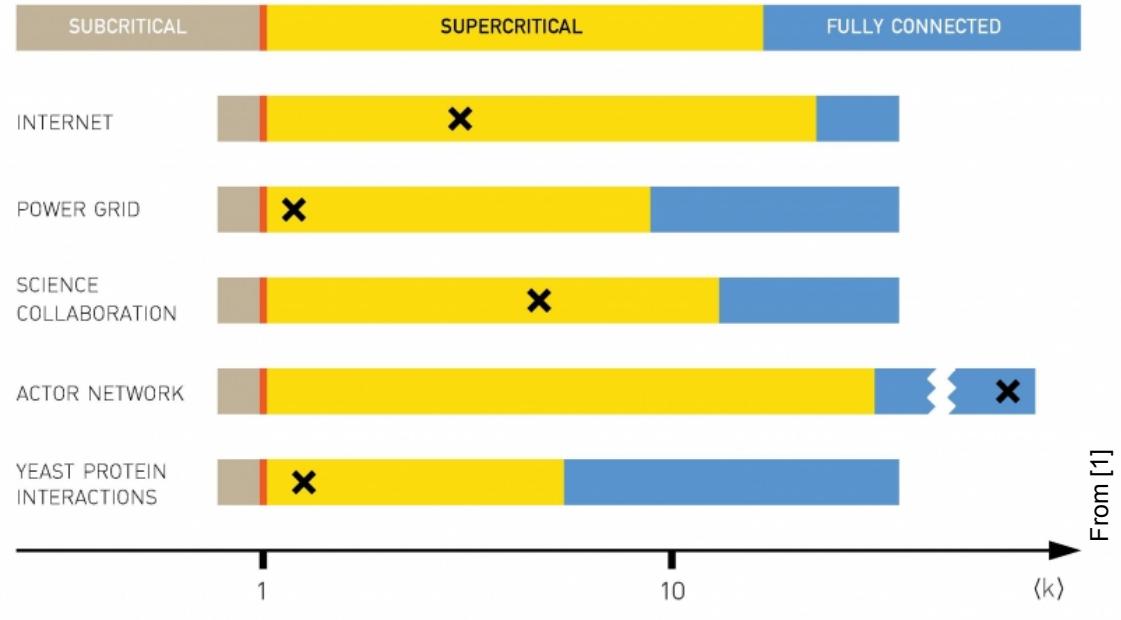
Network	N	L	$\langle k \rangle$	$\ln N$
Internet	192,244	609,066	6.34	12.17
Power Grid	4,941	6,594	2.67	8.51
Science Collaboration	23,133	94,437	8.08	10.05
Actor Network	702,388	29,397,908	83.71	13.46
Protein Interactions	2,018	2,930	2.90	7.61

- In real networks, $\langle k \rangle > 1$: they must have a giant component
- $\langle k \rangle = 1000$ in social networks
- $\langle k \rangle = 7000$ synapses/neuron in brain network

- To form connected networks, $\langle k \rangle > \ln(N)$: most real networks are not connected (according to model)

Most real networks are indeed in the supercritical regime!

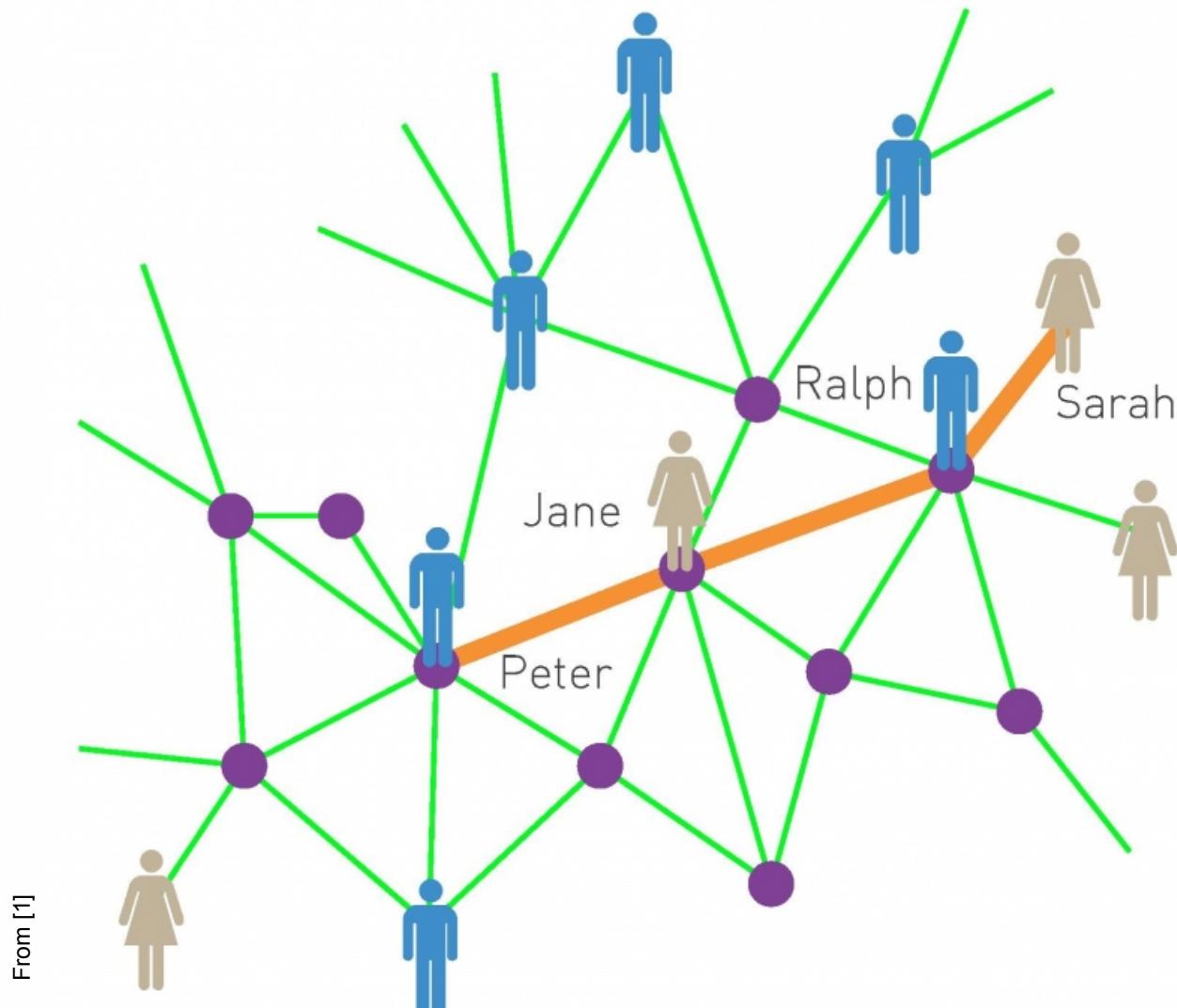
!! Prediction valid under random network assumption !!



Outline

- Why network models?
- Random network model
- Random network evolution
- *Small worlds*
- Clustering coefficient

Small world phenomenon

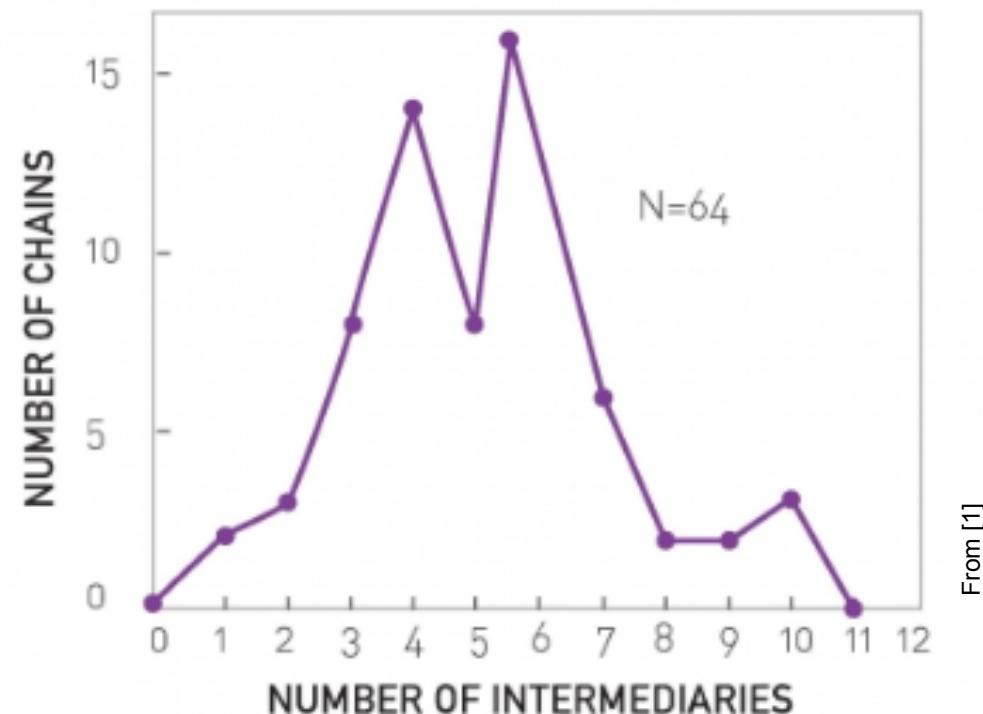


Six degrees of separation:
there is a path of at most six
acquaintances between any
two individuals

Frigyes Karinthy, 1929
Stanley Milgram, 1967

Milgram experiment (1967)

- A stock broker in Boston and a divinity student in Sharon (MA) chosen as targets.
- A letter with information about the target person, is sent to random residents of Wichita and Omaha
- These were asked to forward the letter to a friend, relative or acquaintance who is most likely to know the target person.
- Within a few days the first letter arrived, passing through only two links.
- Eventually 64 of the 296 letters made it back
 - the median number of intermediates was 5.2 (a relatively small number)



Distance between nodes

- Random network with average degree $\langle k \rangle$
- A node has on average $\langle k \rangle$ nodes at distance 1, $\langle k \rangle^2$ at distance 2, $\langle k \rangle^3$ nodes at distance 3, etc...
- The expected number of nodes up to distance d from our starting node is

$$N(d) \approx 1 + \langle k \rangle + \langle k \rangle^2 + \dots + \langle k \rangle^d = \frac{\langle k \rangle^{d+1} - 1}{\langle k \rangle - 1}$$

- This is bounded by the diameter of the network: $N(d_{\max}) \approx N$
- This means that $\frac{\langle k \rangle^{d_{\max}}}{\langle k \rangle \gg 1} \approx N$ or that $d_{\max} \approx \frac{\ln N}{\ln \langle k \rangle}$
The world is indeed small!
 - the largest path is actually dominated by extreme paths, so we rather have the *average distance* as
$$\langle d \rangle \approx \frac{\ln N}{\ln \langle k \rangle}$$
Denser networks get even smaller
 - distances are orders of magnitude smaller than the size of the network (*small* means logarithmic dependence on the size of the network)

Small World Corrections

- The $\langle k \rangle^d$ expansion used to compute diameter must break down as it approaches N , as there are not enough nodes to continue the expansion
- The diameter is better approximated by

$$d_{\max} = \frac{\ln N}{\ln \langle k \rangle} + \frac{2 \ln N}{\ln [-W(\langle k \rangle \exp - \langle k \rangle)]} \quad W \text{ is the Lambert W-function}$$

- The second term is the *correction*, as the number of nodes grows slower when getting closer to the network's diameter
- If $\langle k \rangle \rightarrow 1$, we can calculate

$$d_{\max} = 3 \frac{\ln N}{\ln \langle k \rangle} \quad \text{Increase due to tree-like structure}$$

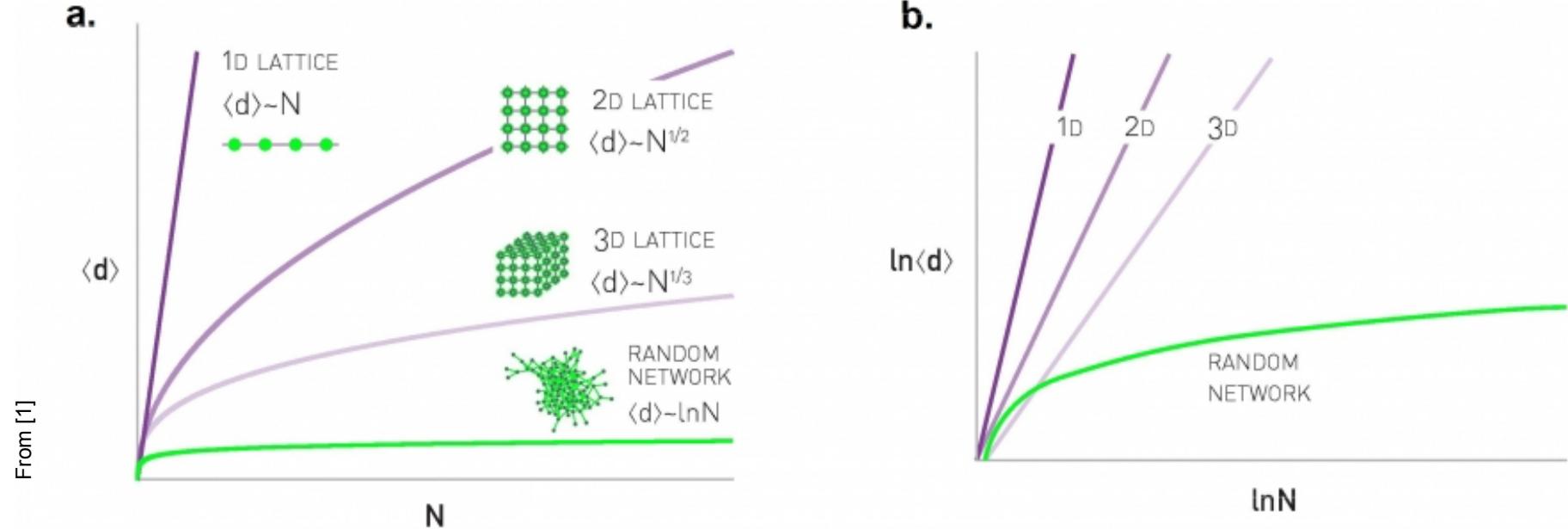
- If $\langle k \rangle \rightarrow \infty$, we have

$$d_{\max} = \frac{\ln N}{\ln \langle k \rangle} + \frac{2 \ln N}{\langle k \rangle} + \ln N \left(\frac{\ln \langle k \rangle}{\langle k \rangle^2} \right)$$

Increasing $\langle k \rangle$ gets closer to initial prediction

D. Fernholz and V. Ramachandran. The diameter of sparse random graphs. Random Structures and Algorithms, 31:482-516, 2007.

Why are small worlds surprising?



- Our intuition about distance is ‘wrongly’ based on regular lattices
 - if a social network would live on a square lattice the average distance between two persons would be $(7 \times 10^9)^{1/2} = 83.666$!!

Distances in reference networks

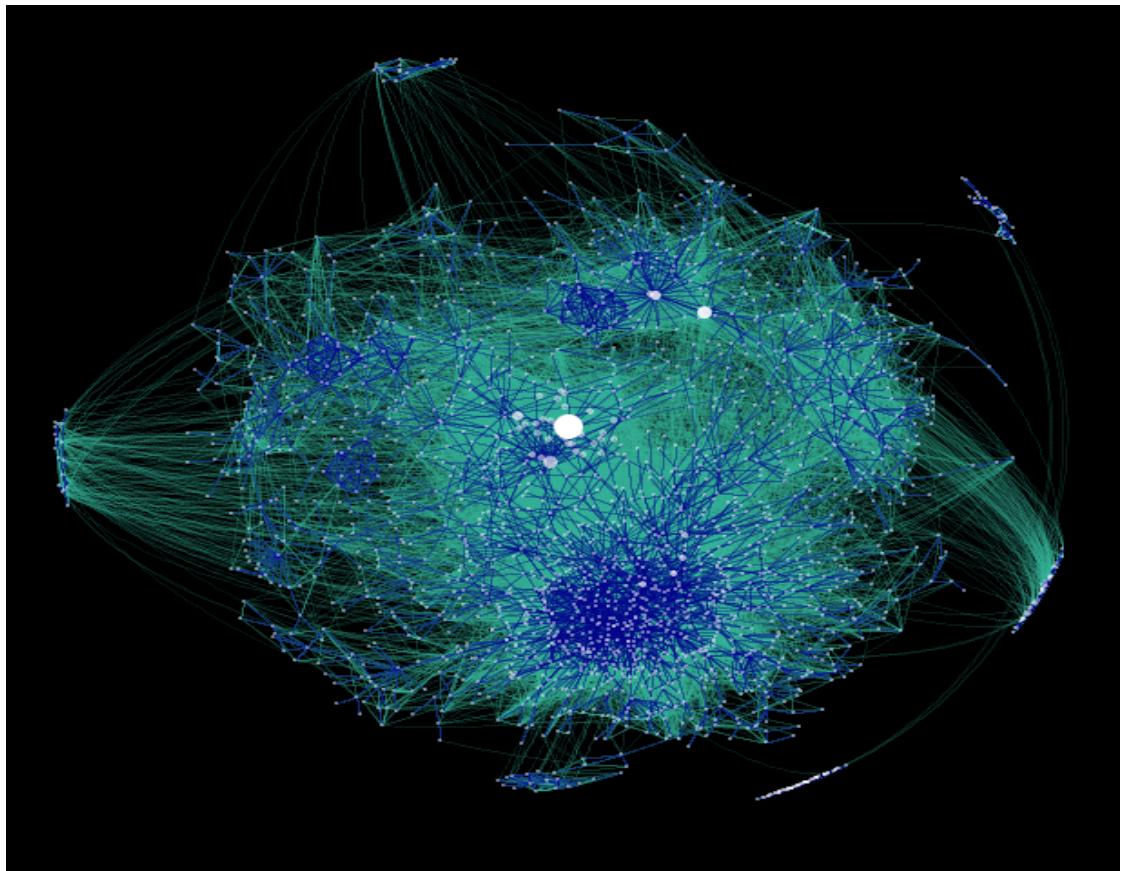
Network	N	L	$\langle k \rangle$	$\langle d \rangle$	d_{\max}	Our prediction $\ln N / \ln \langle k \rangle$
Internet	192,244	609,066	6.34	6.98	26	6.58
WWW	325,729	1,497,134	4.60	11.27	93	8.31
Power Grid	4,941	6,594	2.67	18.99	46	8.66
Mobile-Phone Calls	36,595	91,826	2.51	11.72	39	11.42
Email	57,194	103,731	1.81	5.88	18	18.4
Science Collaboration	23,133	93,437	8.08	5.35	15	4.81
Actor Network	702,388	29,397,908	83.71	3.91	14	3.04
Citation Network	449,673	4,707,958	10.43	11.21	42	5.55
E. Coli Metabolism	1,039	5,802	5.58	2.98	8	4.04
Protein Interactions	2,018	2,930	2.90	5.61	14	7.14

From [1]

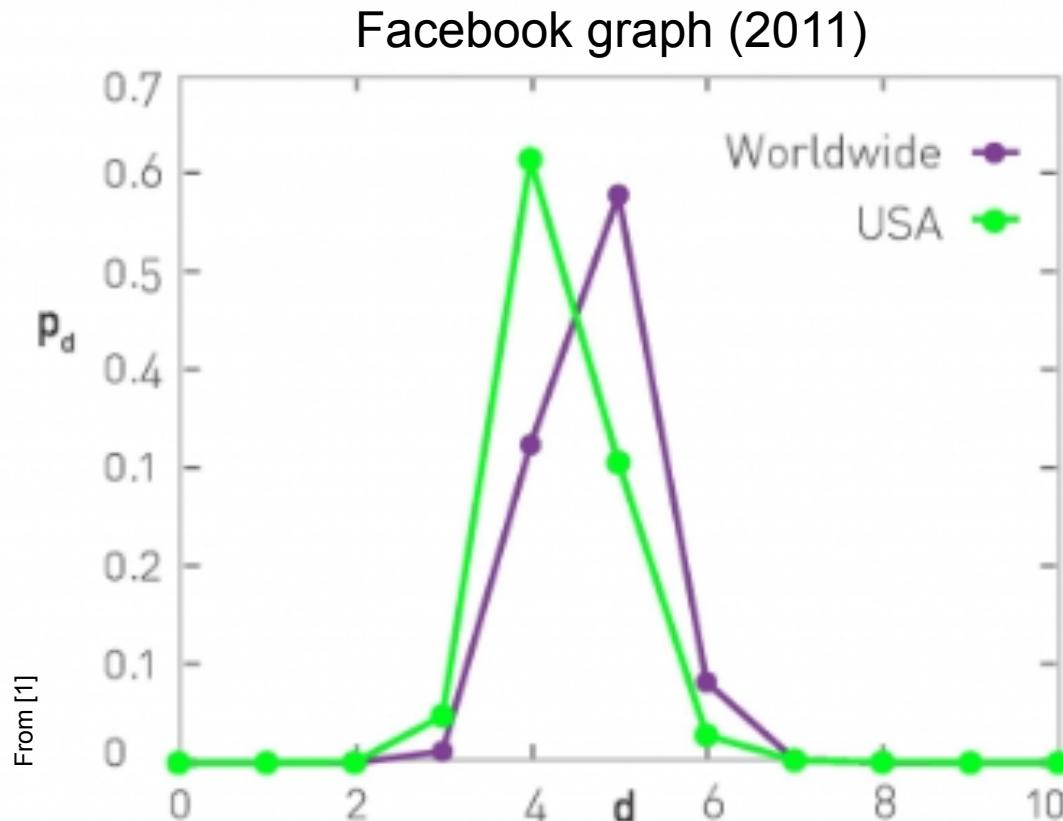
Given the huge differences in scope, size, and average degree, the agreement is excellent!

WWW example

- Number of clicks to reach a random webpage?
 - lack of complete map of web
 - *finite size scaling*: measuring average path length in samples of increasing sizes
$$\langle d \rangle \approx 0.35 + 0.89 \ln N$$
 - in 1999, 800 millions docs: prediction of 19 degrees of separations
$$\langle d \rangle \approx 18.69$$
 - then, measurements on 200 millions docs, give $\langle d \rangle \approx 16$, while prediction is
$$\langle d \rangle \approx 17$$
 - currently, $N \sim 10^{12}$, and estimation would be
$$\langle d \rangle \approx 25$$



Facebook small world



$$N = 7.2 \times 10^8$$

$$L = 6.8 \times 10^{10}$$

$$\langle d \rangle = 4.74$$

- For the global social network:

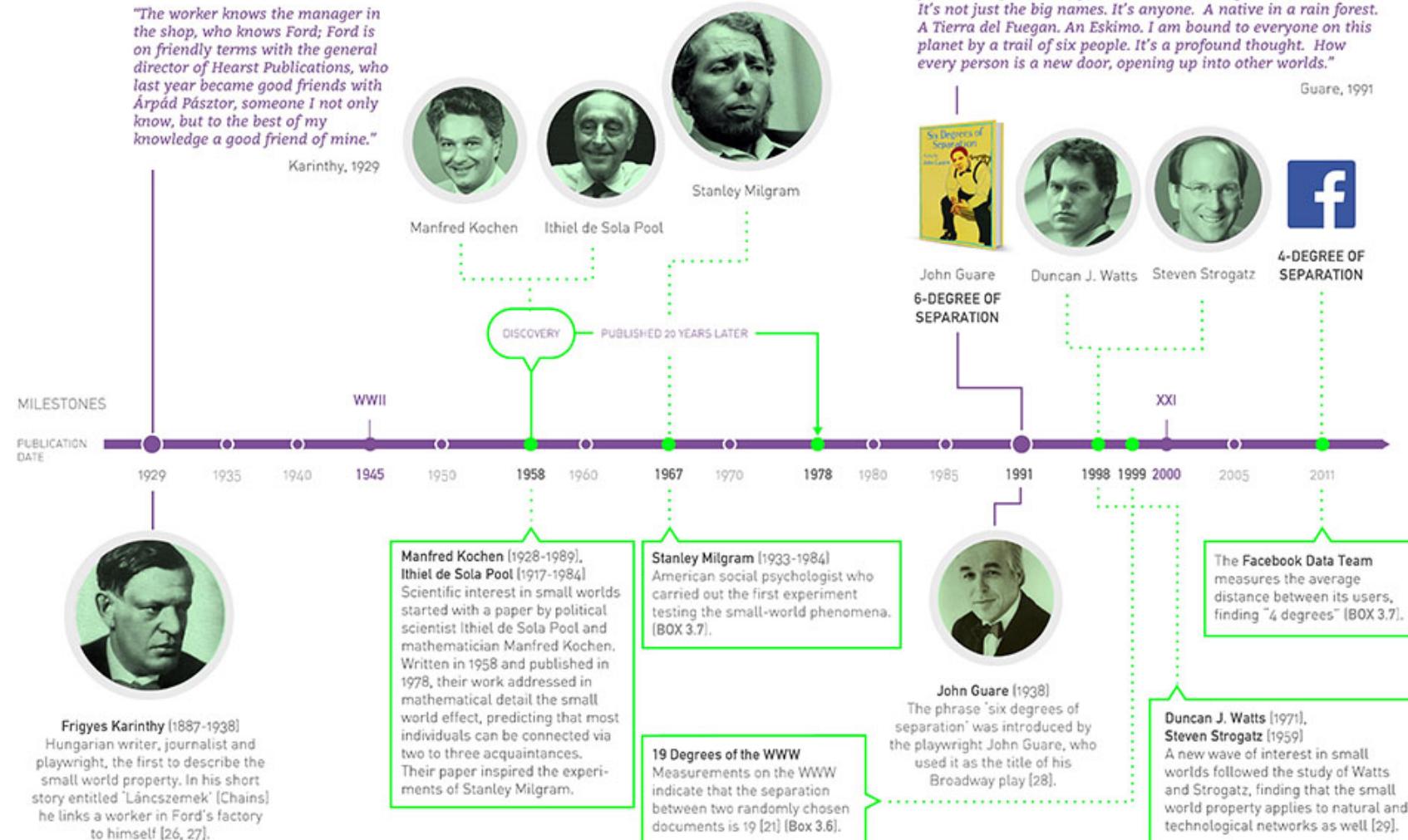
$$\langle k \rangle = 10^3$$

$$N = 7 \times 10^9$$

$$\langle d \rangle \approx \frac{\ln 7 \times 10^9}{\ln(10^3)} = 3.28$$

Closer to reality than Milligram's 6

Small world examples



Outline

- Why network models?
- Random network model
- Random network evolution
- Small worlds
- *Clustering coefficient*

Clustering coefficient

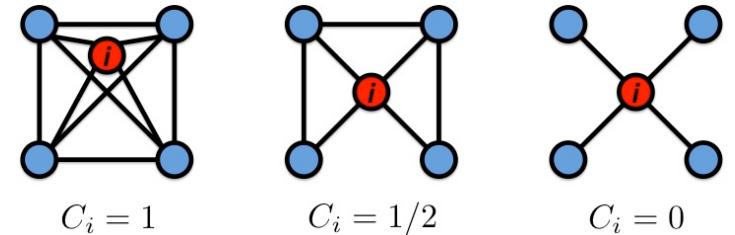
- Number of links between neighbours in random network

$$\langle L_i \rangle = p \frac{k_i(k_i - 1)}{2}$$

- Local clustering coefficient (density of links in the neighbourhood):

$$C_i = \frac{2 \langle L_i \rangle}{k_i(k_i - 1)} = p = \frac{\langle k \rangle}{N}$$

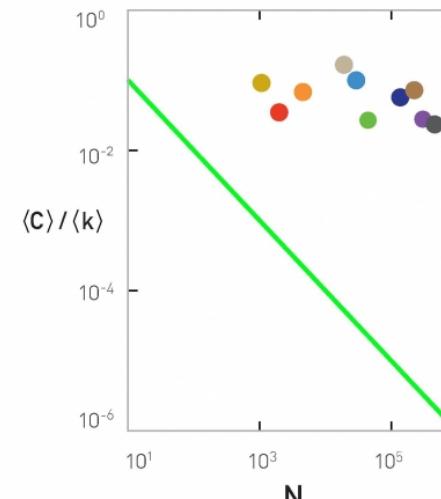
- for fixed $\langle k \rangle$, the larger the network, the smaller the local clustering coefficient, as well as the average one $\langle C \rangle$
- C_i is independent of the node's degree in random networks
- the clustering coefficient is generally small for random networks



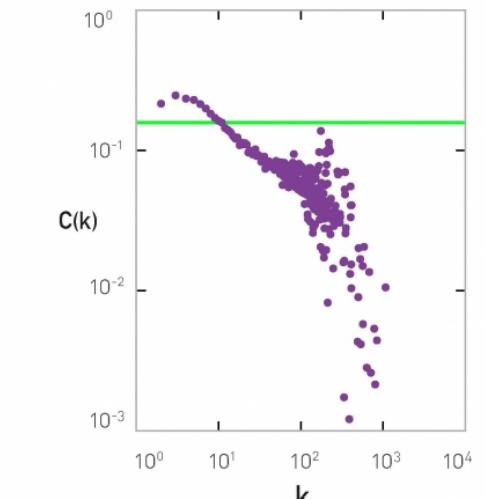
Clustering coefficient in real networks

- In contrary to prediction, $\frac{\langle C \rangle}{\langle k \rangle}$ does not decrease as N^{-1}
 - it rather seems independent of N
- In contrary to prediction, $C(k)$ further decreases with the degree
- The random network model does not capture the clustering of real networks!
 - real networks have a much higher clustering coefficient than predicted

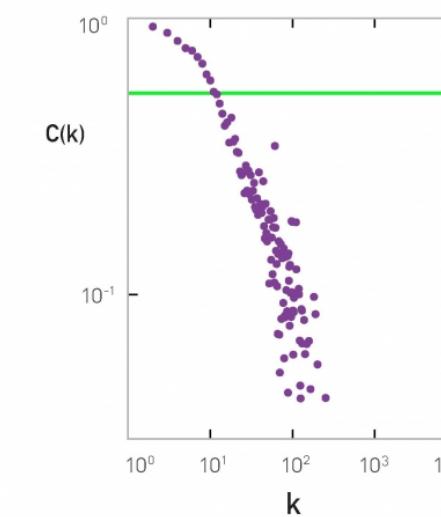
a. All Networks



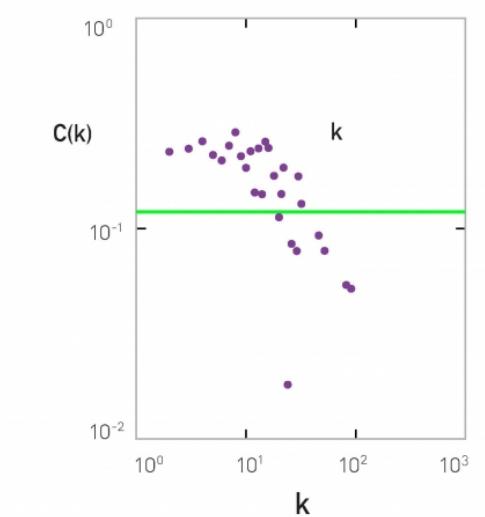
b. Internet



c. Science Collaboration

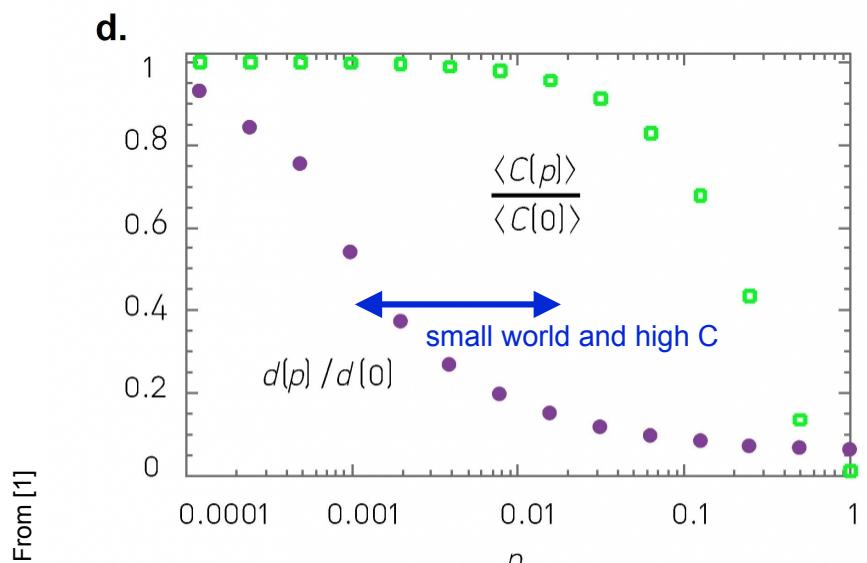
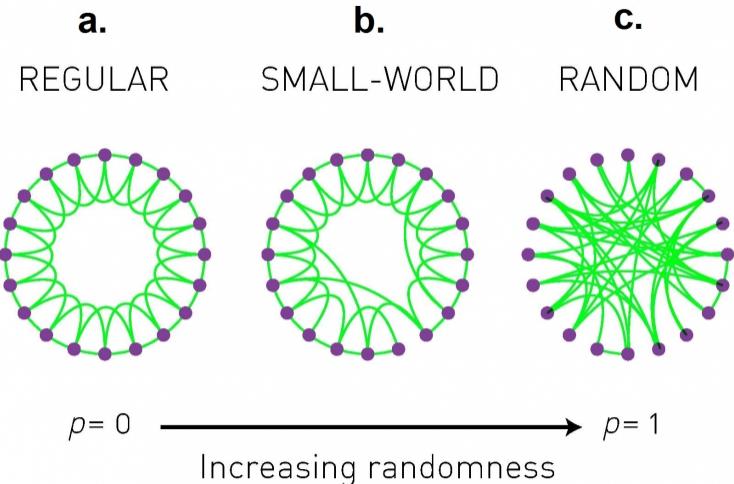


d. Protein Interactions



From [1]

Watts-Strogatz Model



(Normalized wrt regular lattice, and with $N = 1000$ and $\langle k \rangle = 10$)

- Extension of the random network model for
 1. small world property
 2. high clustering
- This *small world* model interpolates between a regular lattice and a random network

- Model construction
 1. start from a ring, $C(k) = 3/4$
 2. rewire each link with probability p
 3. for $p = 1$ all links are rewired into a random network

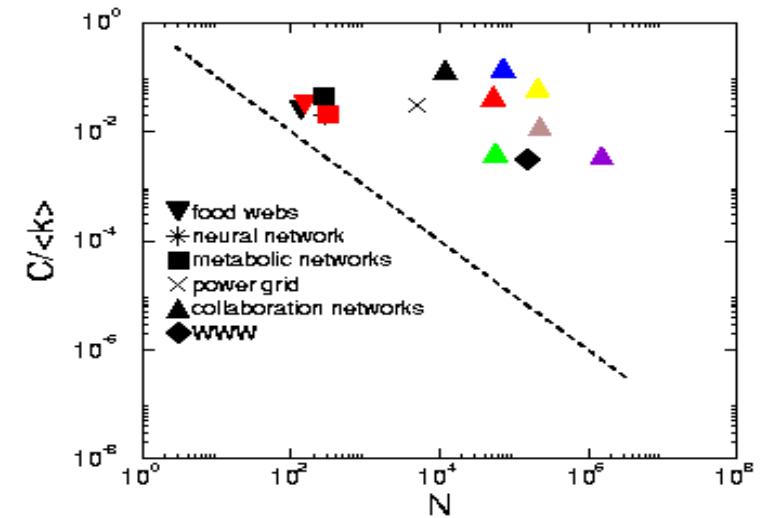
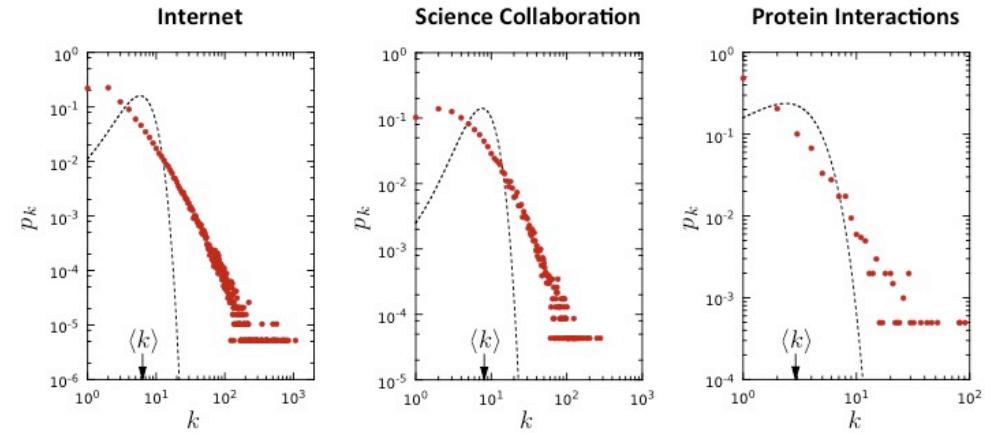
Short path lengths and high clustering coexist in the network

Watts and Strogatz, Nature, 1998

Real networks are not random!

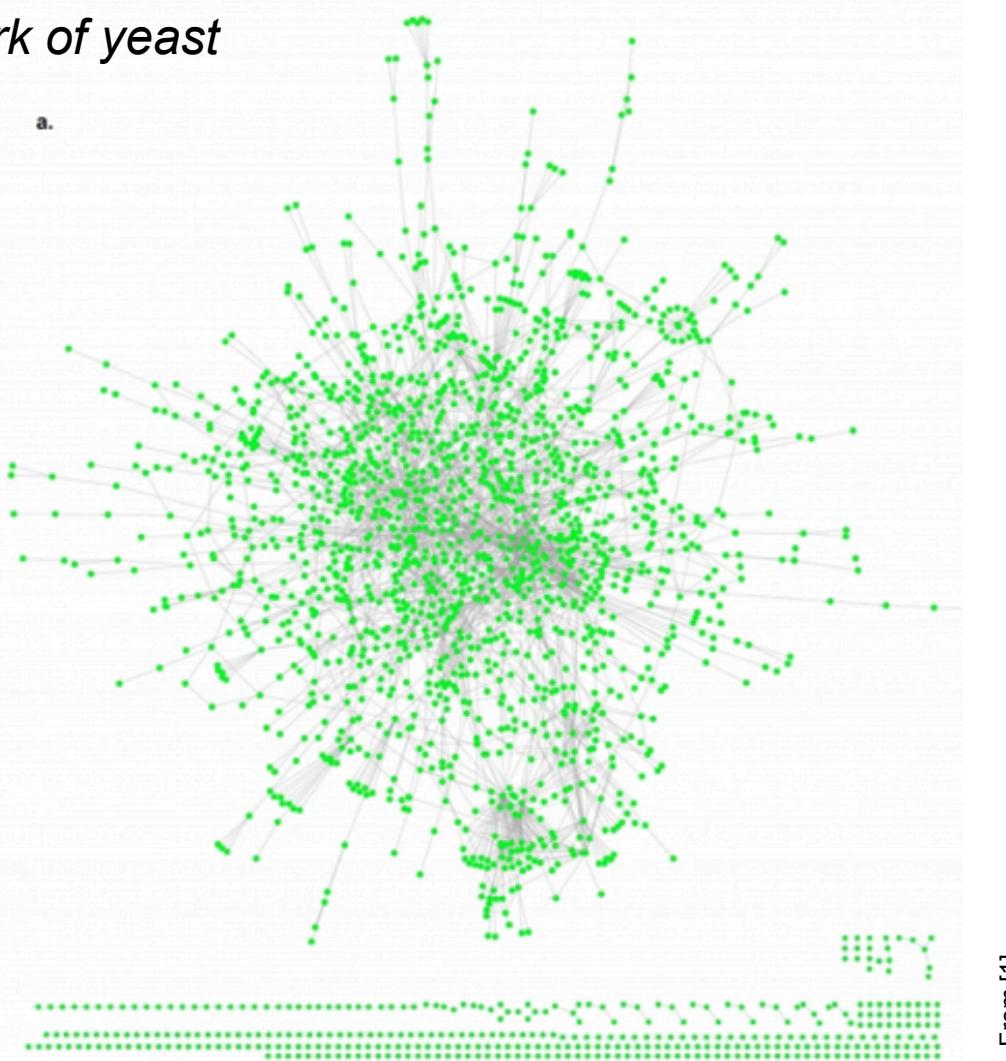
- ✓ Average path length is reasonably well captured by the model (*small world*)

- ◆ Actual degree distribution is not Poisson
 - larger number of highly connected nodes
- ◆ Actual clustering coefficient decreases with degree and is independent of system size
- ◆ Most networks are not broken into isolated clusters even if they indeed have a giant component



This is not a random network?

protein interaction network of yeast



Summary

- Random network: N nodes, with each pair connected with probability p

- Average degree

$$\langle k \rangle = p(N - 1)$$
$$\langle L \rangle = \frac{pN(N - 1)}{2}$$

- Degree distribution

- Binomial form

$$p_k = \binom{N - 1}{k} p^k (1 - p)^{N - 1 - k}$$

- Poisson approximation

$$p_k = e^{-\langle k \rangle} \frac{\langle k \rangle^k}{k!}$$

- Giant component size

$$\begin{cases} \langle k \rangle < 1 : & N_G \sim \ln N \\ 1 < \langle k \rangle < \ln N : & N_G \sim N^{\frac{2}{3}} \\ \langle k \rangle > \ln N : & N_G \sim (p - p_c)N \end{cases}$$

- Average distance

$$\langle d \rangle \propto \frac{\ln N}{\ln \langle k \rangle}$$
$$\langle C \rangle = \frac{\langle k \rangle}{N}$$

- Clustering coefficient

References

- [1] Network Science, by Albert-László Barabási, 2016
- [2] Networks: An Introduction, by M. Newman, 2010

