

# Network Analysis:

The Hidden Structures behind the Webs We Weave

17-213 / 17-668

## Small-World Networks

Tuesday, October 31, 2023

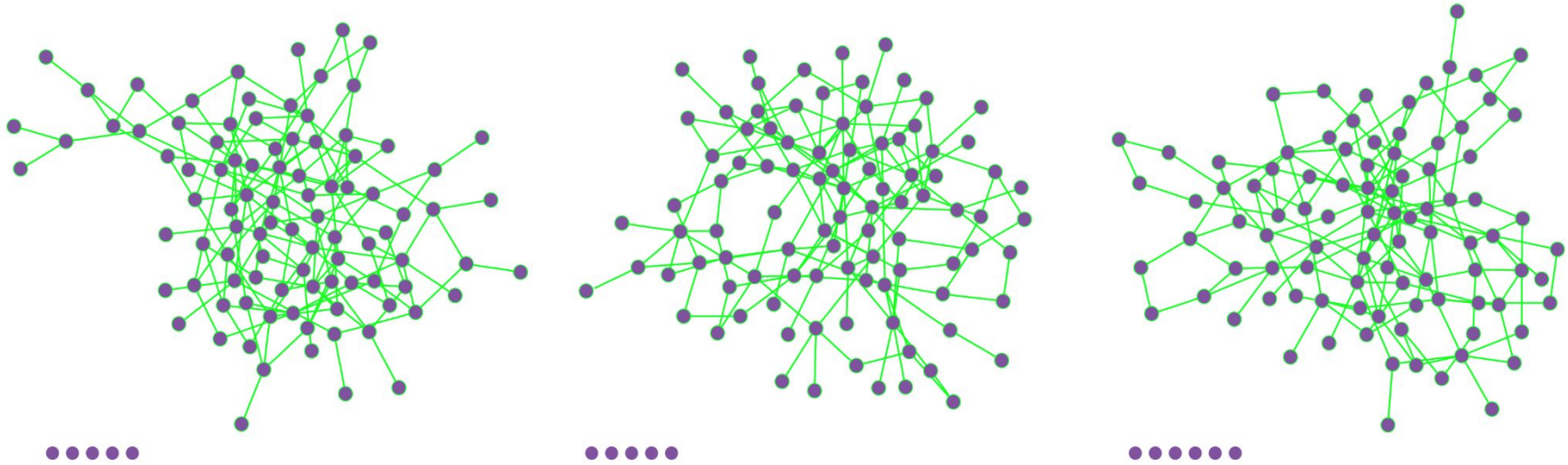
Patrick Park & Bogdan Vasilescu

# 2-min Quiz, on Canvas

# Network Models

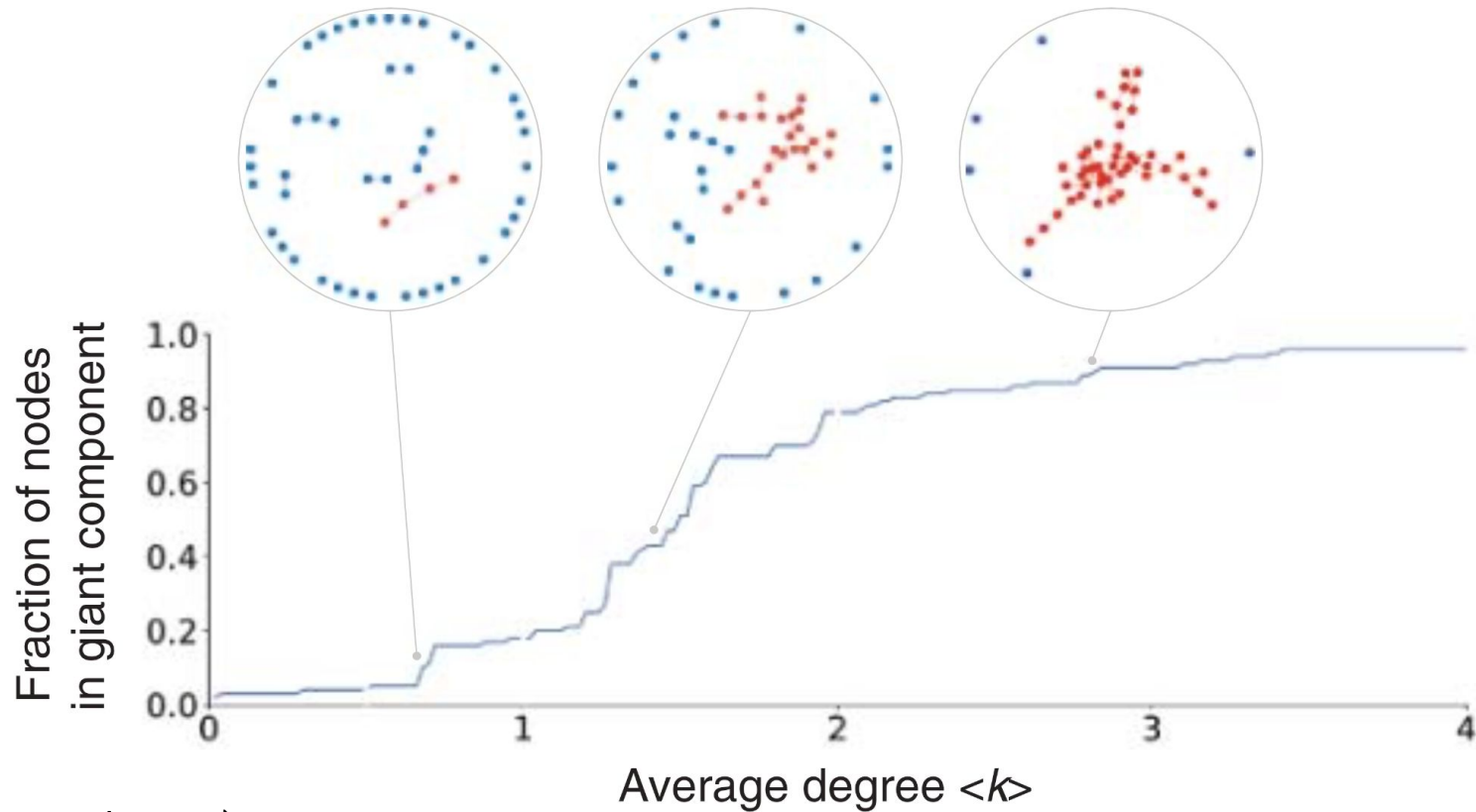
# 1. Random Networks (recall)

A random network consists of  $N$  nodes where each node pair is connected with probability  $p$

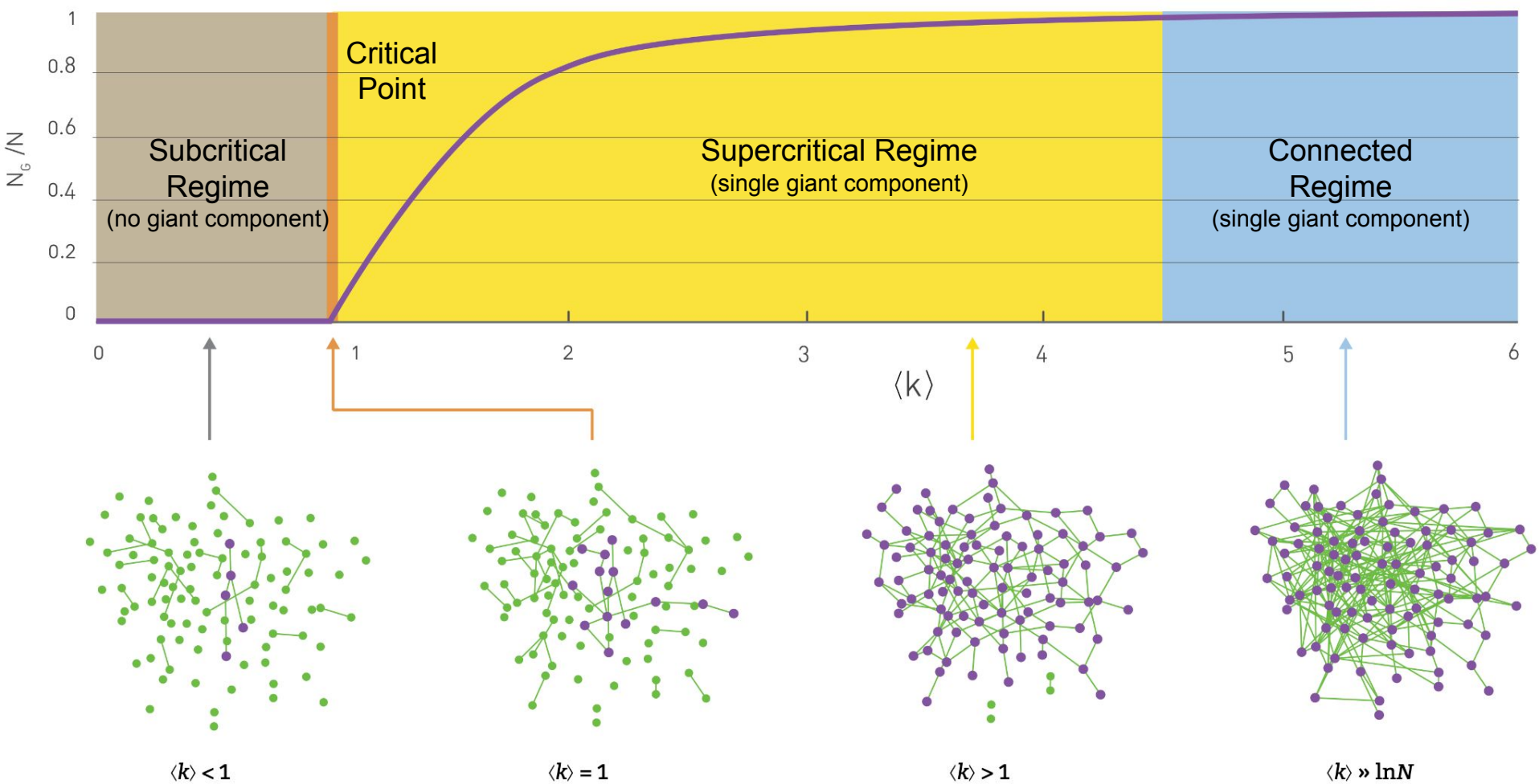


Three realizations of a random network with  $p=0.03$  and  $N=100$ . Several nodes have degree  $k=0$ , shown as isolated nodes at the bottom.

# Around $\langle k \rangle = 1$ a giant component grows very fast



(Menczer et al, 2020)



(Barabasi Ch. 3.6; Erdős & Rényi, 1959 )

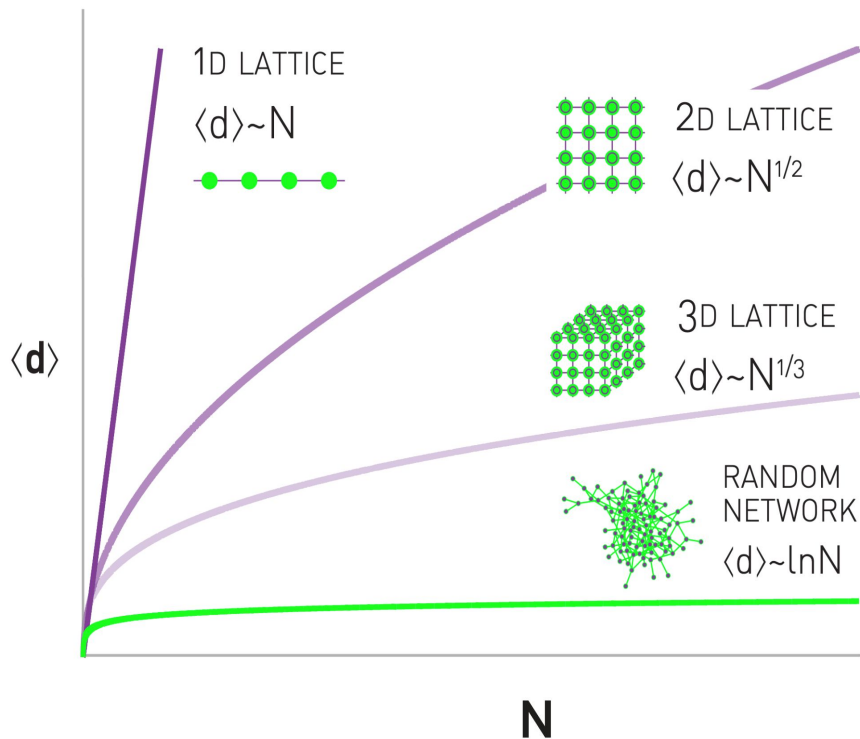
# Random networks are “small worlds”

The dependence of the average distance in a random network on  $N$  and  $\langle k \rangle$ :

$$\langle d \rangle \approx \frac{\ln N}{\ln \langle k \rangle}$$

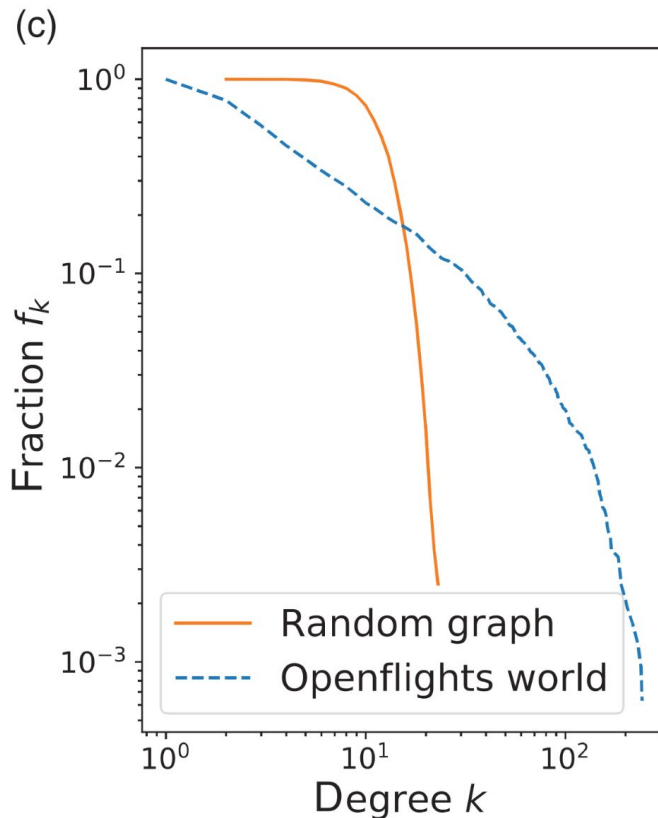
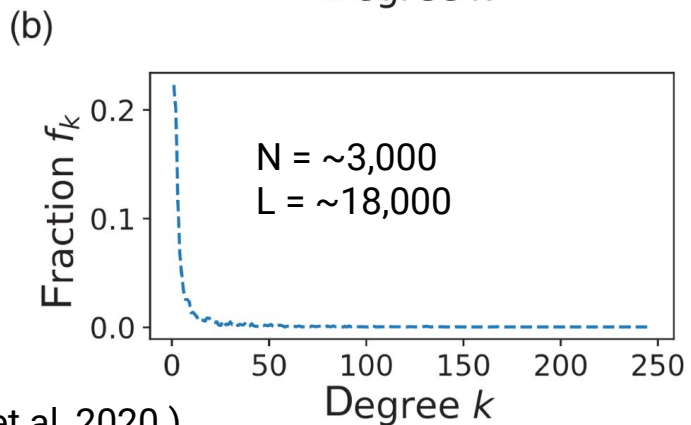
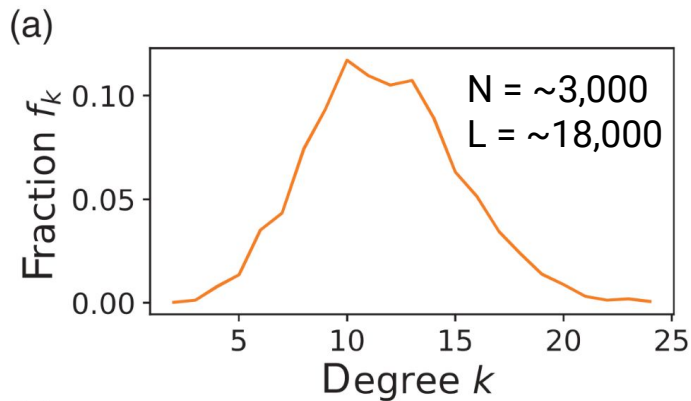
The distances in a random network are orders of magnitude smaller than the size of the network.

(For our world social network, if  $N \approx 7 \times 10^9$  and  $\langle k \rangle \approx 10^3$ , we get  $\langle d \rangle \approx 3.28$ .)

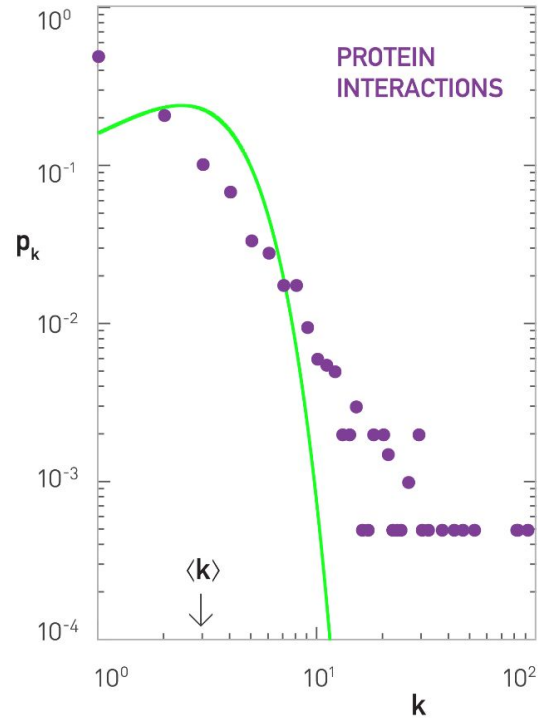
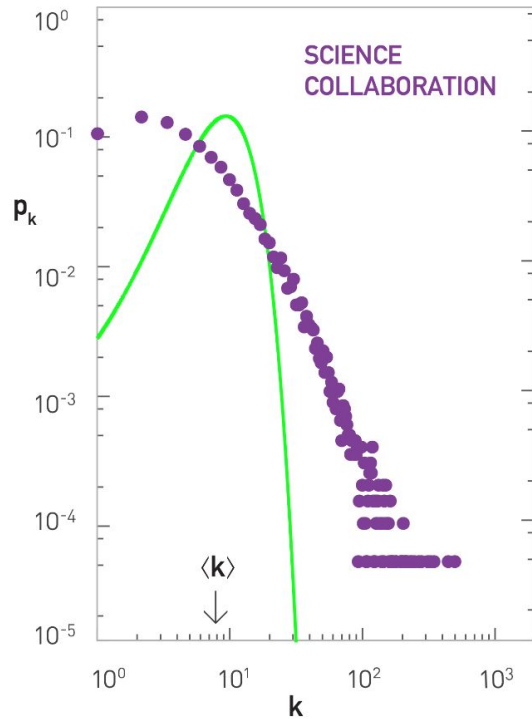
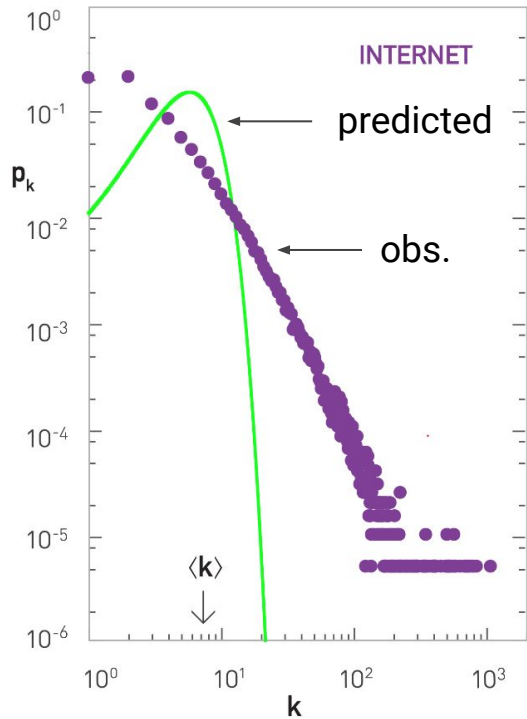




# The random network model underestimates the size and frequency of the high degree nodes, and the number of low degree nodes.



# The random network model underestimates the size and frequency of the high degree nodes, and the number of low degree nodes.



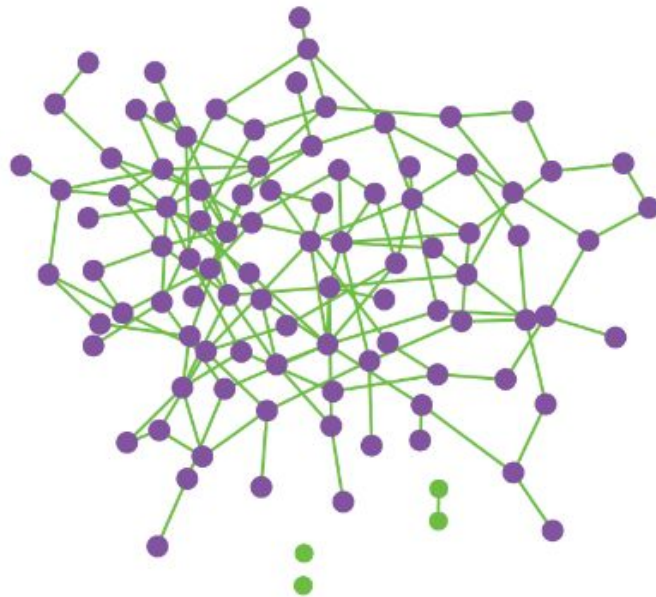
# Probability of a triangle is close to 0 in random networks

Random graphs are useful as a baseline model.

But real-world social networks differ from random graphs in an important way.

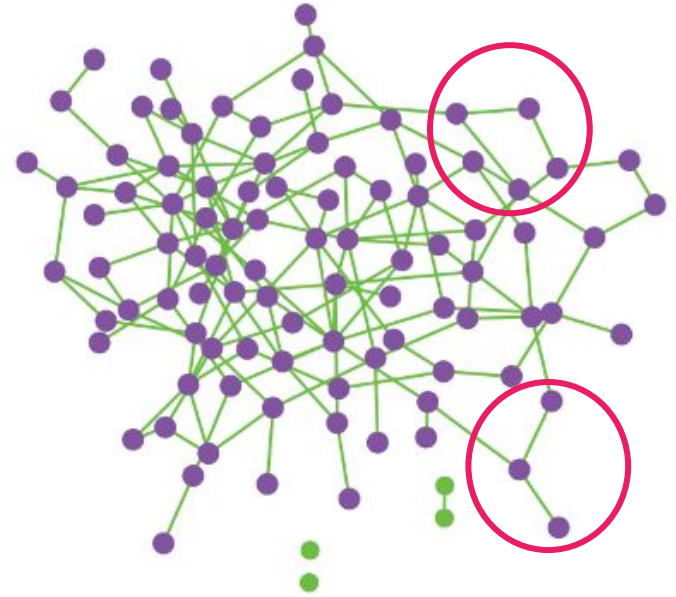
**They contain more triangles.**

These triangles, or triads, are the telltale sign of social groups.



# Aside: Here's why

On a random network, the probability that a pair of neighbors of a node is connected is ...?

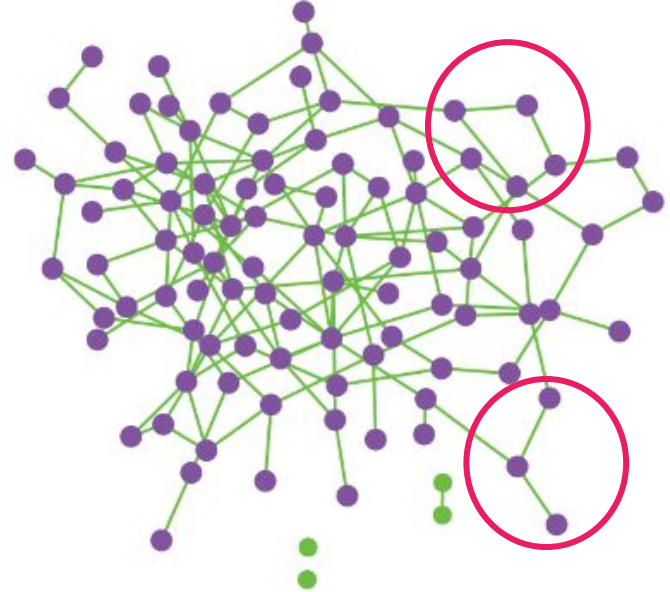


# Aside: Here's why

On a random network, the probability that a pair of neighbors of a node is connected is  $p$ .

The link probability is the same for every pair of nodes by construction ( $p$ ), regardless of their having common neighbors or not.

Thus, the average clustering coefficient is well approximated by  $p$ , a very small number in real networks (because real networks are sparse).



# Summary

## Random Networks

Real networks are different from random ones.

E-R networks do have short paths, but triangles are rare, resulting in average clustering coefficient values that may be orders of magnitude smaller than those measured in real networks.

## 2. Small Worlds

# Short paths & high clustering?

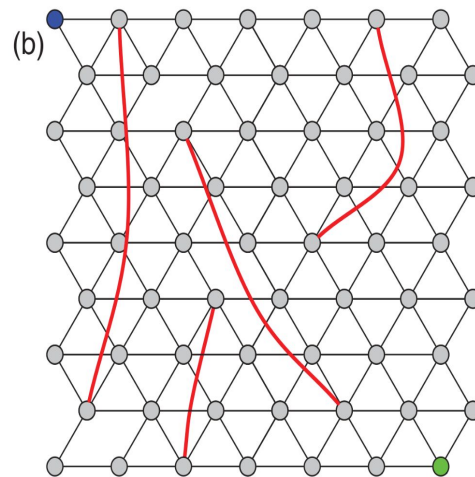
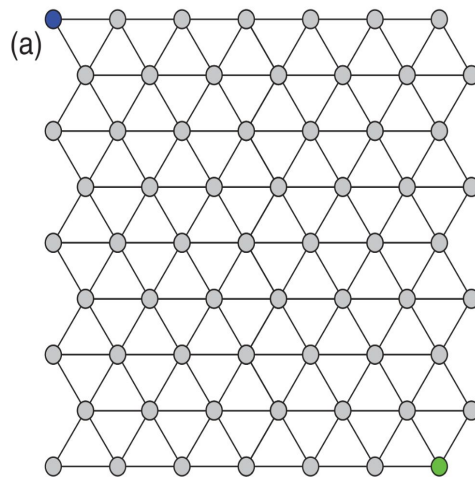
Duncan Watts & Steve Strogatz idea:

Start from a grid-like network where all nodes have the same number of neighbors  $\rightarrow$  high average clustering coefficient.

The internal nodes have degree  $k = 6$  and clustering coefficient  $C = 6/\binom{6}{2} = 6/15 = 2/5$ . Border nodes have smaller degree  $k = 4, 3, 2$  and even higher clustering coefficients, respectively  $C = 3/\binom{4}{2} = 1/2$ ,  $C = 2/\binom{3}{2} = 2/3$ , and  $C = 1/\binom{2}{1} = 1$ . Therefore the average clustering coefficient is at least  $2/5$  and converges to  $2/5$  in the limit of infinite lattice ( $N \rightarrow \infty$ ).

(a) Average path length is large.

(Menczer et al, 2020; Ch. 5)





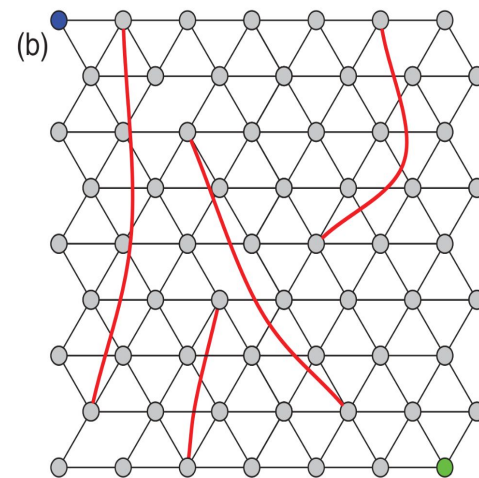
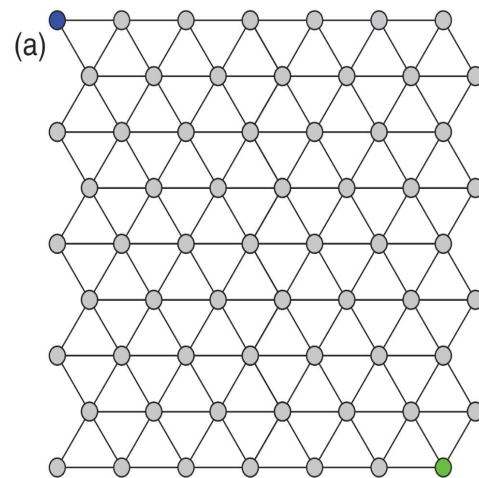
# Short paths & high clustering?

Duncan Watts & Steve Strogatz idea:

(b) Rewire a few links to randomly selected nodes, creating shortcuts (in red).

The shortest path from blue to green goes down substantially.

But the average clustering coefficient remains high, because only a few triangles are disrupted by the rewiring.



# Rewiring procedure: Preserve one endpoint of a randomly chosen link, replace the other endpoint with a node chosen at random

Applies to each link of the network with rewiring probability  $p$ .

The expected number of rewired links is  $pL$ .

$p = 0 \rightarrow$  initial lattice

$p = 1 \rightarrow$  random network

(All links are rewired to random nodes, which is equivalent to placing links between randomly chosen pairs of nodes)

# Rewiring procedure: Preserve one endpoint of a randomly chosen link, replace the other endpoint with a node chosen at random

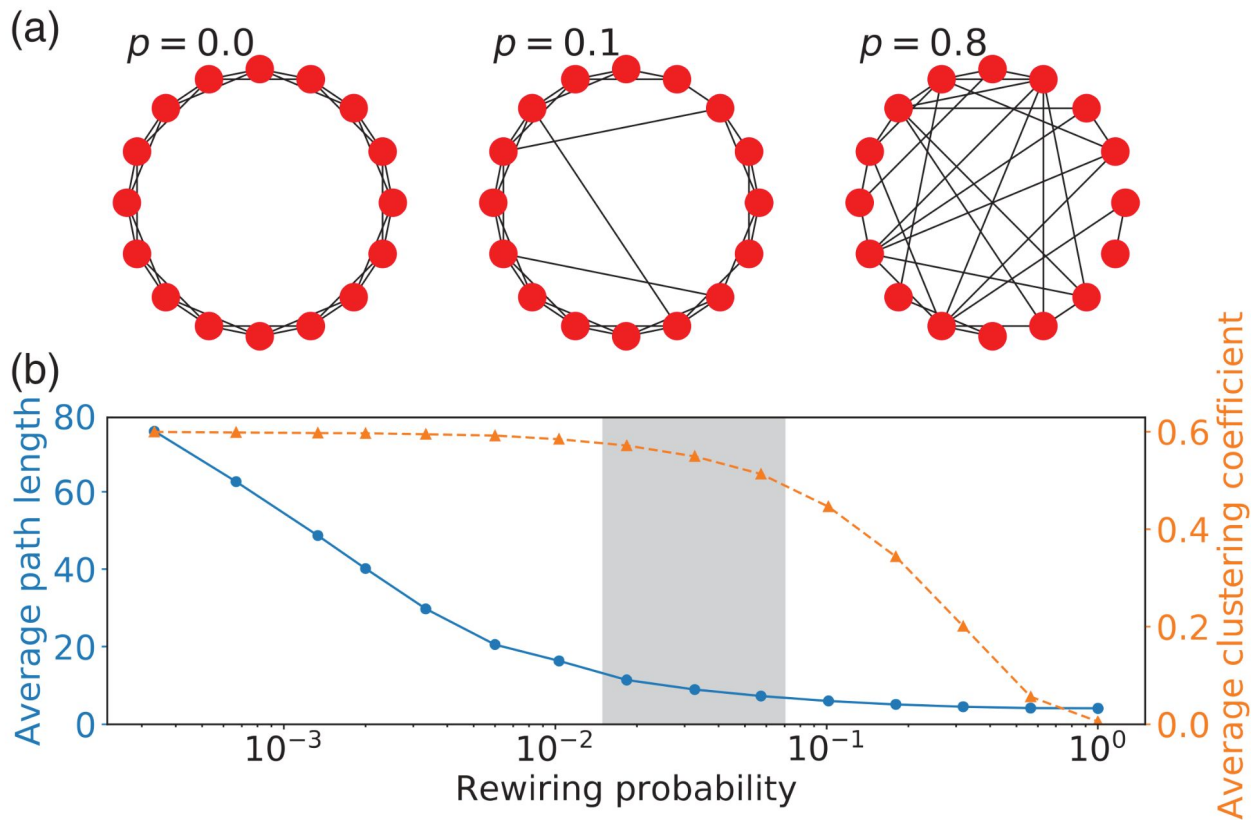
Applies to each link of the network with rewiring probability  $p$ .

The expected number of rewired links is  $pL$ .

$p = 0 \rightarrow$  initial lattice (path length is high)

$p = 1 \rightarrow$  random network (clustering is low)

Sweet spot is grey area: Average path length is almost as short as in random network, clustering coefficient is almost as high as lattice.



# Small World Network Simulation

<https://ccl.northwestern.edu/netlogo/models/SmallWorlds>

# Example: American Corporate Elite (Board Interlock)

**Table 2** Changes in the elite network, 1982–1999

	N	N (component)	K (avg degree)	L (avg geodesic)	C	L (random)	C (random)	SW quotient
All 'large' firms								
1982 boards	648	581	10.0	3.38	0.24	2.76	0.017	11.34
1990 boards	591	524	8.8	3.46	0.24	2.88	0.017	11.87
1999 boards	600	516	8.6	3.46	0.22	2.93	0.016	11.84
1982 directors	6505	5853	19.0	4.27	0.88	2.94	0.003	186.82
1990 directors	5393	4768	17.0	4.30	0.87	2.99	0.004	169.21
1999 directors	5311	4538	16.0	4.33	0.87	3.06	0.003	183.03
Single panel of firms at three points in time								
1982 boards	195	177	6.8	3.15	0.24	2.70	0.039	5.33
1990 boards	195	185	7.6	3.06	0.23	2.58	0.041	4.73
1999 boards	195	186	7.2	2.98	0.20	2.64	0.039	4.55
1982 directors	2366	2179	19.1	4.03	0.91	2.61	0.009	67.23
1990 directors	2078	1976	17.4	3.98	0.89	2.65	0.009	67.26
1999 directors	1916	1819	16.3	3.86	0.88	2.69	0.009	68.35

# Example: German Company Ownership Network

**Table 2. How Small Is Germany's Small World: A Comparison**

Network	Path Length		Clustering		Actual-to-Random Ratio for:		
	Actual	Random	Actual	Random	Length	Clustering	Length/ Clustering
Film actors network <sup>a</sup>	3.65	2.99	.79	.001	1.22	2,925.93	2,396.90
Power grid network <sup>a</sup>	18.70	12.40	.08	.005	1.51	16.00	10.61
<i>C. Elegans</i> network <sup>a</sup>	2.65	2.25	.28	.05	1.18	5.60	4.75
German firms, connected	5.64	3.01	.84	.022	1.87	38.18	22.46
German owners, connected	6.09	5.16	.83	.008	1.18	118.57	100.48

<sup>a</sup> Data come from Watts and Strogatz (1998). See text above for descriptions of these networks.

# Summary

## Small Worlds

### **Still no hubs!**

The degree distribution transitions from that of the initial lattice (all nodes have identical degree), to that of a random network with the same number of nodes and links.

Hence, for any value of the rewiring probability  $p$ , all nodes have similar degree.



**Recall: A scale-free network is a network whose degree distribution follows a power law**

$$\log p_k \sim -\gamma \log k$$

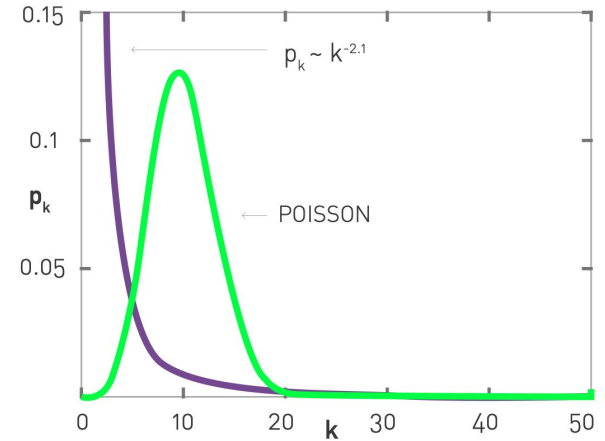
# Poisson vs. Power-law Distributions

Small  $k$ : power law is above the Poisson  $\rightarrow$  a scale-free network has a **large number of small degree nodes**, most of which are absent in a random network.

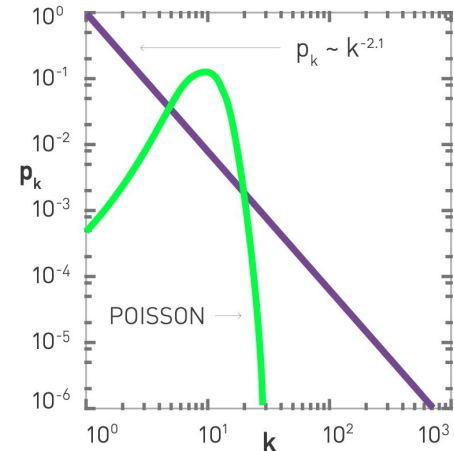
$k$  around  $\langle k \rangle$ : the Poisson is above the power law  $\rightarrow$  in a random network there is an **excess of nodes with degree  $k \approx \langle k \rangle$**

Large  $k$ : power law is again above the Poisson  $\rightarrow$  **observing a high-degree node, or hub, is orders of magnitude more likely in a scale-free network.**

(a)

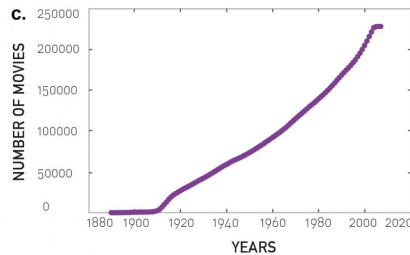
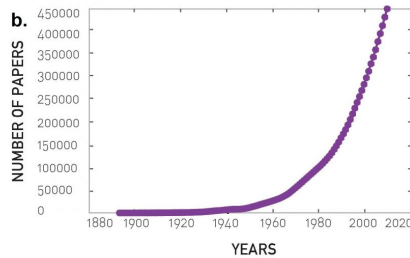
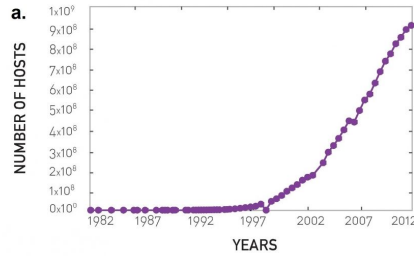


(b)



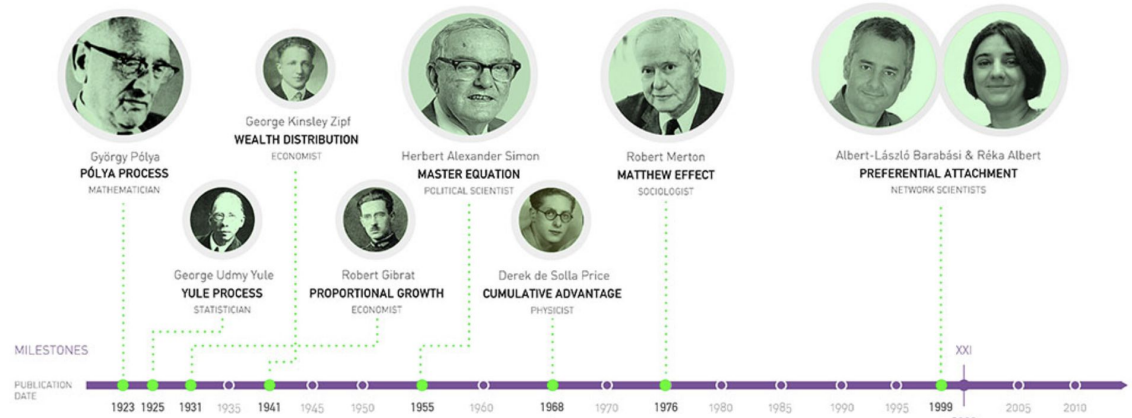
# 3. Preferential Attachment

# The Preferential Attachment (BA) Model Explains Scale-Free



Two assumptions:

- The network grows, one node added at a time
- A new node is more likely to link to high degree nodes
  - Rich get richer, “Matthew effect”, Zipf’s law...



# Preferential Attachment Mechanics

We start from a complete graph with  $m_0$  nodes. Each iteration of the algorithm consists of two steps:

1. A new node  $i$  is added to the network, with  $m \leq m_0$  new links attached to it. The parameter  $m$  is thus the average degree of the network.
2. Each new link is wired to an old node  $j$  with probability

$$\Pi(i \leftrightarrow j) = \frac{k_j}{\sum_l k_l}. \quad (5.9)$$

The denominator in Eq. (5.9) is the sum of the degrees of all nodes (except  $i$ ), and guarantees that the sum of all probabilities equals one, as it must be.

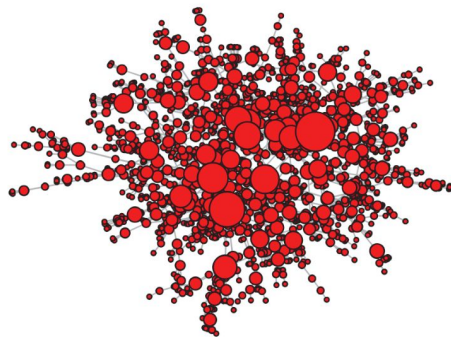
The procedure is repeated until the network reaches the desired number of nodes  $N$ .

# Preferential attachment results in hubs

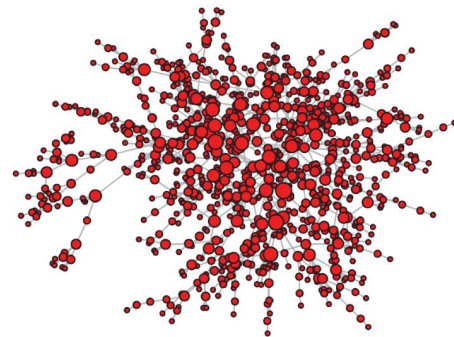
Older nodes get more chances to receive links, which makes them even more likely to attract new links in the future.

Growth alone is not sufficient! (b)

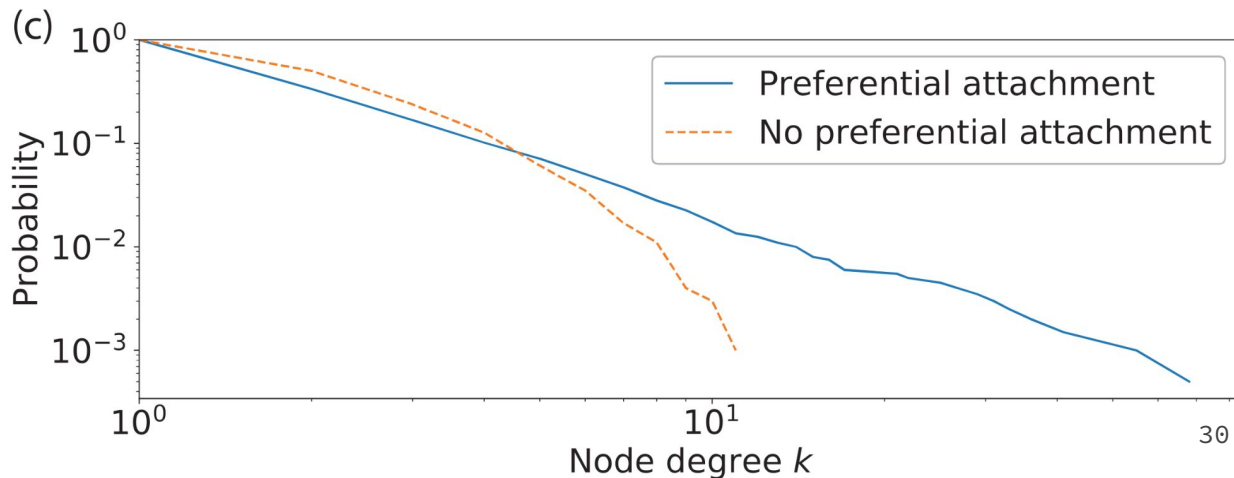
(a)



(b)



(c)



# Non-linear preferential attachment?

$$\Pi_{\alpha}(i \leftrightarrow j) = \frac{k_j^{\alpha}}{\sum_l k_l^{\alpha}}, \quad (5.10)$$

where the exponent  $\alpha$  is a parameter. For  $\alpha = 1$  we recover the standard BA model. What happens when  $\alpha \neq 1$ ? There are two different scenarios:

1. If  $\alpha < 1$ , the link probability does not grow with degree as fast as in the BA model, so the advantage of high-degree nodes over the others is not as big. As a result, the degree distribution does not have a heavy tail — the hubs disappear!
2. If  $\alpha > 1$ , high-degree nodes accumulate new links much faster than low-degree nodes. As a consequence, one of the nodes will end up being connected to a fraction of all other nodes. The effect is even more extreme when  $\alpha > 2$ , in which case we observe a *winner-takes-all* effect: a single node may be connected to all other nodes, which have approximately the same, low degree.

# Linear preferential attachment is the way to go!

$$\Pi_{\alpha}(i \leftrightarrow j) = \frac{k_j^{\alpha}}{\sum_l k_l^{\alpha}}, \quad (5.10)$$

where the exponent  $\alpha$  is a parameter. For  $\alpha = 1$  we recover the standard BA model. What happens when  $\alpha \neq 1$ ? There are two different scenarios:

1. If  $\alpha < 1$ , the link probability does not grow with degree as fast as in the BA model, so the advantage of high-degree nodes over the others is not as big. As a result, the degree distribution does not have a heavy tail — the hubs disappear!
2. If  $\alpha > 1$ , high-degree nodes accumulate new links much faster than low-degree nodes. As a consequence, one of the nodes will end up being connected to a fraction of all other nodes. The effect is even more extreme when  $\alpha > 2$ , in which case we observe a *winner-takes-all* effect: a single node may be connected to all other nodes, which have approximately the same, low degree.



# Summary

Preferential Attachment (BA)

Fixed pattern for the degree distribution: same slope for any choice of model parameters.

→ Real degree distributions could decay faster or more slowly.

# Summary

Preferential Attachment (BA)

The hubs are the oldest nodes.

→ New nodes cannot overcome them in degree.

# Summary

## Preferential Attachment (BA)

It does not create many triangles.

→ The average clustering coefficient is much lower than in many real networks.

# Summary

Preferential Attachment (BA)

Nodes and links are only added.

→ In real networks they can also be deleted.

# Summary

## Preferential Attachment (BA)

Since each node is attached to older nodes, the network consists of a single component.

→ Many real networks have multiple components.

# Other Preferential Models

## 4. Attractiveness Model

BA: If a node has no neighbors, it will never have neighbors!

Idea: Besides degree, make nodes receive links also because of an intrinsic *attractiveness*.

The attractiveness model is a slightly modified version of the original BA model, in which Eq. (5.9) expressing the probability that an old node  $j$  receives a link from a new node  $i$  is replaced by

$$\Pi(i \leftrightarrow j) = \frac{A + k_j}{\sum_l (A + k_l)}, \quad (5.11)$$

where  $A$  is the attractiveness parameter and can take any positive value. The case  $A = 0$  yields the BA model.

# 5. Fitness Model

BA: Hubs are the oldest nodes. But newcomers can overtake existing nodes in popularity. Previous attractiveness parameter is the same for all nodes.

Idea: Model *individual node fitness*.

The fitness model is similar to the BA model, but each node  $i$  is assigned a fitness value  $\eta_i > 0$  generated from some distribution  $\rho(\eta)$ . Then at every step, each new link from a new node  $i$  is wired to an old node  $j$  with probability

$$\Pi(i \leftrightarrow j) = \frac{\eta_j k_j}{\sum_l \eta_l k_l}. \quad (5.12)$$

If all nodes have identical fitness, the model reduces to the BA model, as the constant  $\eta$  is a factor that cancels out between numerator and denominator of Eq. (5.12), returning the standard prescription of preferential attachment.



## 6. Random Walk Model

BA: Triangles are formed rarely, because the probability of a node receiving a link is proportional to its degree, regardless of whether the new pair of neighbors have a common neighbor or not.

Idea: In addition to creating random connections, also connect to a new neighbor's neighbors.

## 6. Random Walk Model

The random walk model can start from any small network. Each iteration of the algorithm consists of the following steps:

1. A new node  $i$  is added to the network, with  $m > 1$  new links attached to it.
2. The first link is wired to an old node  $j$ , chosen at random.
3. Each other link is attached to a randomly selected neighbor of  $j$ , with probability  $p$ , or to another randomly selected node, with probability  $1 - p$ .

The parameter  $p$  is the probability of triadic closure, because by setting a link between  $i$  and a neighbor of  $j$ , say  $l$ , we close the triangle  $(i, j, l)$ . If  $p = 0$ , there is no triadic closure and new nodes choose their neighbors entirely at random. When  $p = 1$  all links except the first one are wired to neighbors of the initially selected old node, thus closing triangles.

# 7. Copy Model

Triadic closure → An individual *copies* the contacts of somebody else.

Copying takes place in other contexts, e.g., gene duplication, citations, ...

The copy model is similar to the previous random walk model (a new node gets wired either to a randomly selected old node, with some probability, or else to its neighbors).

However, there is no triadic closure in the copy model.

# 8. Rank Model

BA: Need to know the degree of the nodes. More realistic to have a perception of the *relative ranking*.

The rank model can start from any small graph with  $m_0$  nodes. A node property, such as the degree, age, or some measure of fitness is selected to rank the nodes. Each iteration of the algorithm consists of the following steps:

1. All nodes are ranked based on the property of interest. Nodes are assigned ranks  $R = 1, 2, \text{etc.}$  Node  $l$  receives rank  $R_l$ .
2. A new node  $i$  is added to the network, with  $m \leq m_0$  new links attached to it.
3. Each new link from  $i$  is wired to an old node  $j$  with probability

$$\Pi(i \leftrightarrow j) = \frac{R_j^{-\alpha}}{\sum_l R_l^{-\alpha}}, \quad (5.13)$$

where the exponent  $\alpha > 0$  is a parameter.

# Summary

# Applications

# Small-World and System-Level Collaboration and Creativity

How does the “small-worldness” (high clustering, short path length) affect a creative industry’s overall collaboration and creativity?

Small-world networks emerge from cohesive structures through **repeated collaborations between artists** (high clustering) and **new collaborations** between artists who belonged to separate clusters (low path length).

Higher “small-worldness” predicts a musical’s success to a certain point

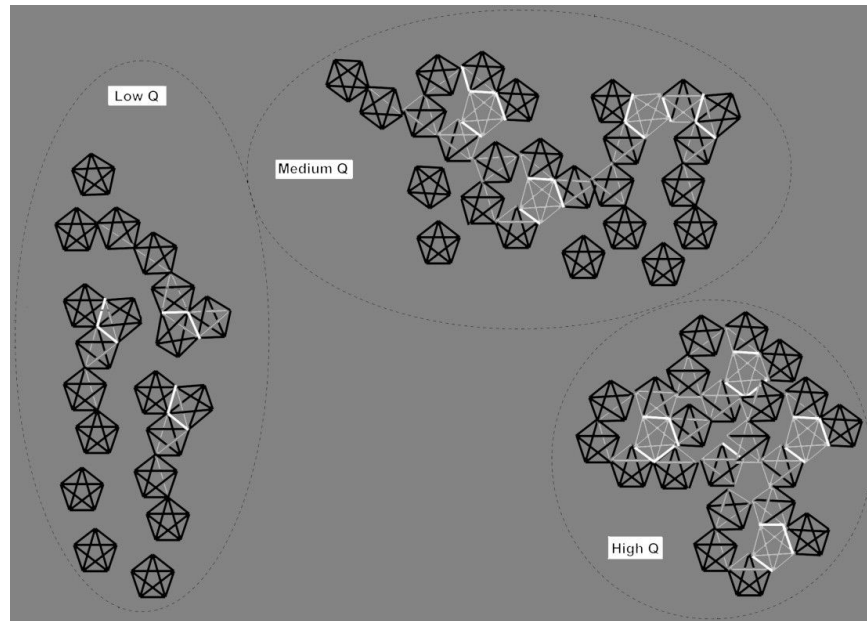
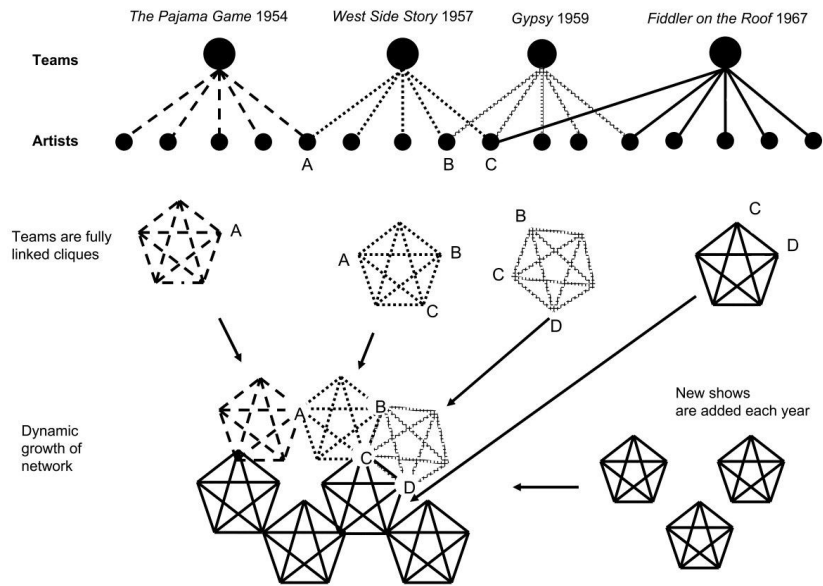
- New ideas can spread efficiently (low PL)
- Risk assessment and implementation cost of those ideas are lowered by familiar collaboration ties (high CC)

Too high “small-worldness” implies low system-level diversity

- New ideas spread too quickly
- Highly cohesive groups fall into “group think”

Prediction: Inverted U-shape relationship between success and “small-worldness”

# Broadway musical artist network





# Measuring “small-worldness”

Small-world quotient  $Q$ :

Ratio of  $CC$  to  $PL$ , where each is normalized by corresponding quantities computed from random bipartite graphs of same size and degree distribution

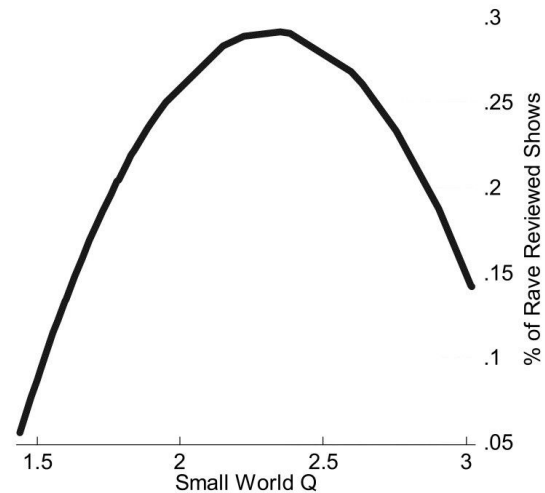
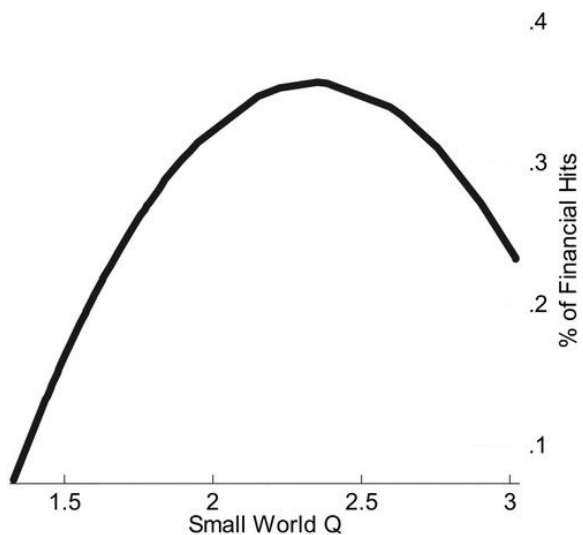
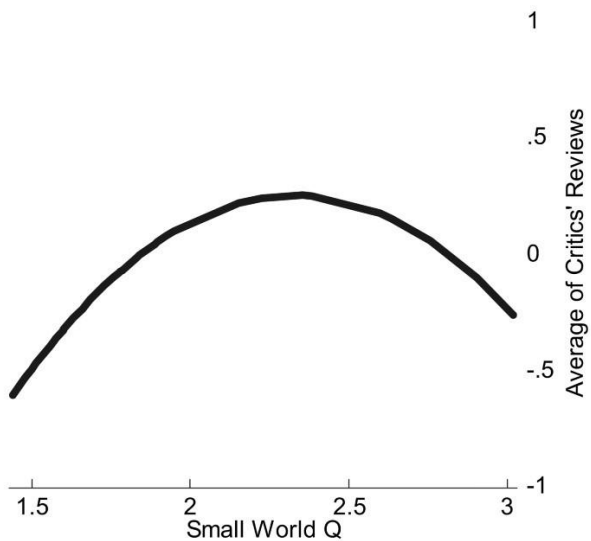
Random bipartite graph generation:

- $M$  musicals' degree distribution (i.e., number of artists)
- $N$  artists' degree distribution (i.e., number of musicals)
- Construct random bipartite graph that holds the two degree distributions constant

# Small-World and System-Level Collaboration and Creativity

YEAR	NEW MUSICALS	AVERAGE TEAM SIZE	CLUSTER COEFFICIENT			PATH LENGTH			$\frac{Q}{CCr/PLr}$	1966 ...	13	7.5	.301	.146	2.05	3.04	2.37	1.28	1.59
			Actual	Random	Ratio	Actual	Random	Ratio											
1945 ...	14	7.1	.287	.077	3.7	3.13	2.23	1.40	2.63	1967 ...	7	8.3	.302	.148	2.04	2.98	2.33	1.27	1.59
1946 ...	19	6.8	.295	.073	4.02	3.11	2.24	1.38	2.90	1968 ...	16	7.9	.329	.165	1.98	2.96	2.32	1.27	1.55
1947 ...	12	6.7	.311	.074	4.15	3.12	2.27	1.37	3.01	1969 ...	13	7.3	.33	.166	1.97	2.97	2.36	1.25	1.57
1948 ...	16	6.6	.315	.078	4.04	3.14	2.34	1.34	3.01	1970 ...	14	7.0	.331	.167	1.98	2.97	2.36	1.25	1.57
1949 ...	9	5.8	.319	.089	3.55	3.04	2.36	1.29	2.75	1971 ...	17	6.0	.354	.18	1.96	3.2	2.46	1.30	1.51
1950 ...	16	8.4	.325	.097	3.34	3.09	2.40	1.28	2.59	1972 ...	16	6.7	.381	.188	2.02	3.53	2.51	1.40	1.43
1951 ...	11	6.4	.33	.109	3.02	3.06	2.41	1.26	2.38	1973 ...	12	7.4	.389	.193	2.01	3.48	2.56	1.36	1.47
1952 ...	8	7.0	.328	.109	3.01	3.03	2.36	1.28	2.35	1974 ...	9	6.4	.391	.189	2.06	3.54	2.57	1.37	1.49
1953 ...	8	8.0	.338	.116	2.9	2.98	2.28	1.30	2.22	1975 ...	17	7.3	.371	.146	2.54	3.74	2.58	1.44	1.75
1954 ...	11	7.9	.328	.115	2.85	2.98	2.24	1.32	2.15	1976 ...	14	7.7	.376	.146	2.57	3.75	2.58	1.45	1.77
1955 ...	12	7.8	.33	.133	2.47	2.93	2.22	1.31	1.88	1977 ...	7	7.0	.375	.139	2.69	3.72	2.53	1.47	1.82
1956 ...	10	7.8	.345	.135	2.55	2.93	2.24	1.31	1.94	1978 ...	19	6.6	.364	.141	2.57	3.61	2.49	1.44	1.78
1957 ...	11	7.2	.355	.14	2.53	2.97	2.25	1.31	1.92	1979 ...	16	8.6	.358	.148	2.41	3.42	2.45	1.39	1.72
1958 ...	11	7.0	.353	.136	2.58	3.06	2.25	1.36	1.89	1980 ...	14	7.9	.365	.149	2.43	3.54	2.49	1.42	1.71
1959 ...	16	7.4	.342	.139	2.45	3.03	2.27	1.33	1.84	1981 ...	17	7.8	.355	.153	2.31	3.48	2.53	1.37	1.68
1960 ...	13	6.8	.343	.144	2.38	3.06	2.34	1.31	1.81	1982 ...	15	7.1	.355	.169	2.09	3.57	2.53	1.41	1.48
1961 ...	17	6.1	.338	.146	2.31	3.09	2.38	1.29	1.78	1983 ...	10	8.6	.361	.178	2.02	3.61	2.59	1.39	1.45
1962 ...	13	6.0	.344	.16	2.13	3.14	2.39	1.31	1.63	1984 ...	4	7.0	.358	.178	2	3.59	2.58	1.39	1.44
1963 ...	13	6.9	.324	.154	2.09	3.21	2.38	1.34	1.55	1985 ...	9	7.9	.366	.183	2	3.58	2.61	1.37	1.46
1964 ...	17	6.9	.314	.147	2.12	3.17	2.38	1.33	1.59	1986 ...	8	7.5	.369	.173	2.12	3.69	2.61	1.41	1.50
1965 ...	18	7.2	.304	.141	2.15	3.12	2.40	1.29	1.65	1987 ...	8	6.3	.383	.207	1.85	3.71	2.67	1.39	1.33
										1988 ...	8	6.9	.409	.2	2.04	3.87	2.69	1.43	1.42
										1989 ...	10	6.9	.406	.182	2.23	3.6	2.62	1.37	1.62

# Financial and artistic success of Broadway shows



# Randomization Strategies

# Randomizing requires care

## Strategy of comparing against random networks (deeper dive)

- [Other randomization strategies](#) for small-world networks
- Randomizing nodes vs. edges
  - The case of homophily (nodes vs. weights)
- Hierarchical modeling strategy
  - Successively more constrained random networks as baseline
  - The case of high school romantic network structure

# Considerations for random graph generation

## Know your graph

One-mode vs. two-mode

Directed vs. undirected

Weighted vs. unweighted

→  $2*2*2 = 8$  types

Table 1. Available and Recommended Measurement for Different Types of Networks.

	Path Length	Clustering	Randomization Procedures	Recommended
1-mode network, undirected, unweighted	Unweighted	Global clustering coefficient	Erdos and Renyi Tie rewiring	Tie rewiring (preserves degree distribution)
1-mode network, directed, unweighted	Unweighted	Global clustering coefficient	Erdos and Renyi Tie rewiring	Tie rewiring (preserves degree distribution)
1-mode network, undirected, weighted	Weighted, normalized weights	Global clustering coefficient	Erdos and Renyi Tie rewiring Weight rewiring	Tie rewiring (preserves degree distribution) or weight rewiring (preserves largest connected component, but do not alter structure)
1-mode network, directed, weighted	Weighted, normalized weights	Global clustering coefficient	Erdos and Renyi Tie rewiring Weight rewiring	Tie rewiring (preserves degree distribution) or weight rewiring (preserves largest connected component, but do not alter structure)
2-mode network, undirected, unweighted	Projection, weighted, normalized weights	2-mode global clustering coefficient	Erdos and Renyi Tie rewiring	Tie rewiring (preserves degree distribution)
2-mode network, directed, unweighted	Projection, weighted, normalized weights	2-mode global clustering coefficient	Erdos and Renyi Tie rewiring	Tie rewiring (preserves degree distribution)
2-mode network, undirected, weighted	Projection, asymmetric weights, normalized weights	2-mode global clustering coefficient	Erdos and Renyi Tie rewiring Weight rewiring	Tie rewiring (preserves degree distribution) or weight rewiring (preserves largest connected component, but do not alter structure)
2-mode network, directed, weighted	Projection, asymmetric weights, normalized weights	2-mode global clustering coefficient	Erdos and Renyi Tie rewiring Weight rewiring	Tie rewiring (preserves degree distribution) or weight rewiring (preserves largest connected component, but do not alter structure)

Table 4. Comparison With Different Type of Random Networks.

	U.S. Power Grid	U.S. Airports	C. Elegans Neural Network	Online Social Network	Scientific Collaboration	Online Forum
Transformed networks						
Observed network distance	18.99	2.99	2.46	3.06	6.63	1.88
Expected network distance	8.66	2.51	2.17	2.82	5.59	1.34
Difference (%)	119	19	13	8	19	40
Network distance in classic random networks	8.48 [8.25; 8.70]	2.77 [2.74; 2.80]	2.46 [2.43; 2.49]	3.09 [3.07; 3.11]	5.79 [5.76; 5.82]	1.82 [1.82; 1.82]
Network distance in tie reshuffled random networks	8.49 [8.42; 8.57]	2.59 [2.57; 2.61]	2.38 [2.36; 2.40]	3.04 [3.03; 3.06]	4.94 [4.94; 4.95]	1.86 [1.85; 1.86]
Observed clustering coefficient	0.10	0.35	0.18	0.06	0.36	0.50
Expected clustering coefficient	0.00	0.02	0.05	0.01	0.00	0.18
Difference (%)	18,992	1,371	293	640	105,578	186
Clustering coefficient in classic random network	0.00 [0.00; 0.00]	0.02 [0.02; 0.03]	0.05 [0.04; 0.05]	0.01 [0.01; 0.01]	0.00 [0.00; 0.00]	0.18 [0.18; 0.18]
Clustering coefficient in tie reshuffled random networks	0.00 [0.00; 0.00]	0.24 [0.23; 0.25]	0.11 [0.11; 0.12]	0.08 [0.08; 0.08]	0.01 [0.01; 0.01]	0.43 [0.43; 0.43]

## Classic random vs. tie reshuffling

Classic: pick two nodes at random and connect them

Reshuffling: Start from observed network, randomly pick two edges and swap

(A-B and C-D becomes A-C and B-D)

# Considerations for random graph generation

## Degree distribution

- Model-based: Uniform random (Erdős–Rényi), poisson, power-law...
- Observation-based:
  - Best-fit parametric distribution (power-law with estimated exponent)
  - Observed degree probability density
  - Degree sequence: Each node's degree is preserved



# Node vs. edge randomization

Vast majority randomizes edges

In some studies, node attributes are randomized while edges are not

- Christakis and Fowler's obesity study randomizes node attribute (obesity) to assess obesity clustering
- Measuring homophily level of multiple groups in a network requires edge randomization, preserving each node's degree
  - Average degree of the group can affect homophily measurement
  - Therefore, average degree of the groups need to be "controlled for" in the random graphs

# Summary

Random network construction  
requires careful assessment