

BIG DATA COURSE

DYNAMIC NETWORK PROCESSES

Marc Santolini and Liubov Tupikina

27 Nov 2019

CRI Paris



CANCELED

Planning			
9	18 Dec (3 hours) - Anirudh, Felix	Individual Digital Drivers Individual Project	Open
10	8 Jan (3 hours) - Anirudh, Felix	2. Working with the digital tools	Working with the digital tools
11	15 Jan (3 hours) - Anirudh, Felix	Project Individual/Group Research Project	tutors and strong and 3 weak (student's project)
12	22 Jan (3 hours) - Anirudh, Felix	Individual/Group Mental Health Research Project	I. Work session - group work continues.
			Final presentations and discussion

OUTLINE: OTHER CHANGES

12
December
→
MOVED
TO
JANUARY

Planning		
N°	Type (CM/TP/TD) & hourly rate	
	Describe here the content of each session. (1. means first part of the course, 2. the second part, ~1h20 each)	
1	16 October (3 hours) - Marc, Liubov, Anirudh, Felix	<p>Introduction to Networks and Big Data Efforts</p> <ol style="list-style-type: none"> 1. Introduction to the network part of the course “Why network science?” with presentation of topics covered 2. Introduction to the big data part, with a focus on big data for mental health and presentation of topics covered
2	23 October (3 hours) - Marc, Liubov	<p>Networks - basics of network analysis and visualisation</p> <p>Network metrics and data analysis & create teams for personal projects</p> <ol style="list-style-type: none"> 1. Theory (network construction, centralities, statistical significance, networkx, other packages) 2. Hands-on session working on a dataframe, creating a network, making some statistical analysis, network visualisation
3	6 Nov (3 hours) - Loic	<p>Big Data - Infrastructure of big data 1</p> <ul style="list-style-type: none"> • Introduction to data engineering, what is the day-to-day jobs and what skills are needed
4	13 Nov (3 hours) - Loic	<p>Big Data - Infrastructure of big data 2</p>



Data engineering

Planning		
5	20 Nov (3 hours) - Marc, Liubov, Raphael	<p>Networks - mobility, web-based data</p> <ol style="list-style-type: none"> 1. Hands-on Liuba mobility / geographically embedded networks 2. Hands-on How to get data from Youtube/Twitter/etc APIs for data analysis (Raphael Tackx)
6	27 Nov (3 hours) - Marc, Liubov	<p>Networks - spreading processes</p> <ol style="list-style-type: none"> 1. Dynamics of networks: network growth, network attack, hands on session 2. Dynamics on networks: information spreading (eg spread of fake news, softwares), hands on session
7	4 Dec (3 hours) - Marc, Liubov	<p>Networks - advanced topics</p> <ol style="list-style-type: none"> 1. Temporal networks, networks with attributes: multilayer, multiplex, simplicial, etc. 2. Intro to statistics and data science with R
8	11 Dec (3 hours) - Marc, Liubov	<p>Networks EVALUATION part 1</p> <p>Reverse classroom: 10 minutes presentation of a contribution to Wikipedia about an advanced topic / a paper related to network science. Show slides including modifications of wikipedia page, max 10 slides (1 min per slide)</p>
9	18 Dec (3 hours) - Marc, Liubov	<p>Network EVALUATION part 2</p> <p>10 min project presentation (visualisation and descriptive analysis of network data) through an interactive demonstration of a project notebook including images</p>



TEXTBOOK: CH. 3-5 AND 8

<http://networksciencebook.com/>

Network Science

by Albert-László Barabási

- Personal Introduction
- 1. Introduction
- 2. Graph Theory
- 3. Random Networks
- 4. The Scale-Free Property
- 5. The Barabási-Albert Model

- 6. Evolving Networks
- 7. Degree Correlations
- 8. Network Robustness
- 9. Communities
- 10. Spreading Phenomena
- Preface

TODAY'S GOAL

- 1.Understand how networks can change in time
2. Exercises to handle dynamic processes on networks

b

HOW TO GROW A NETWORK

FROM RANDOM TO PREFERENTIAL ATTACHMENT

After duplication

Proteins



Before duplication

Genes

Genes

FROM RANDOM TO PREFERENTIAL ATTACHMENT

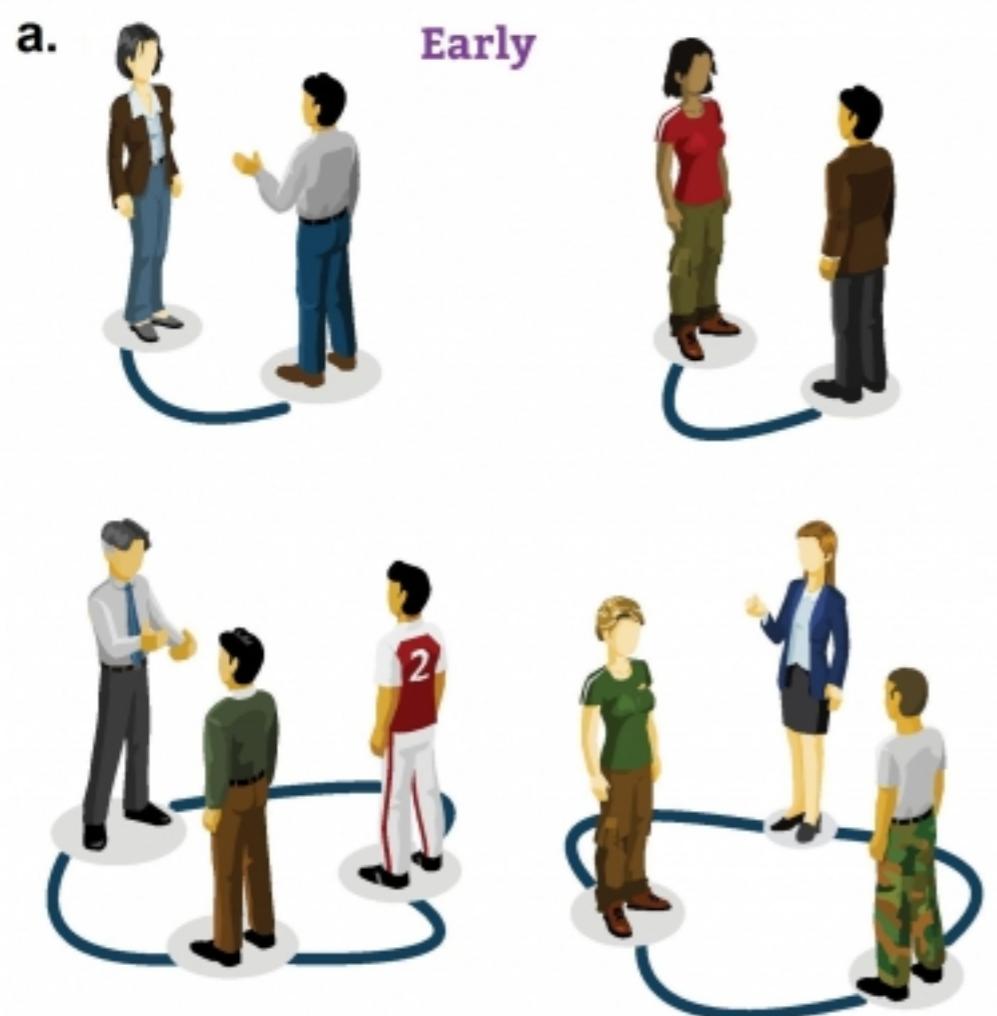


b

After duplication

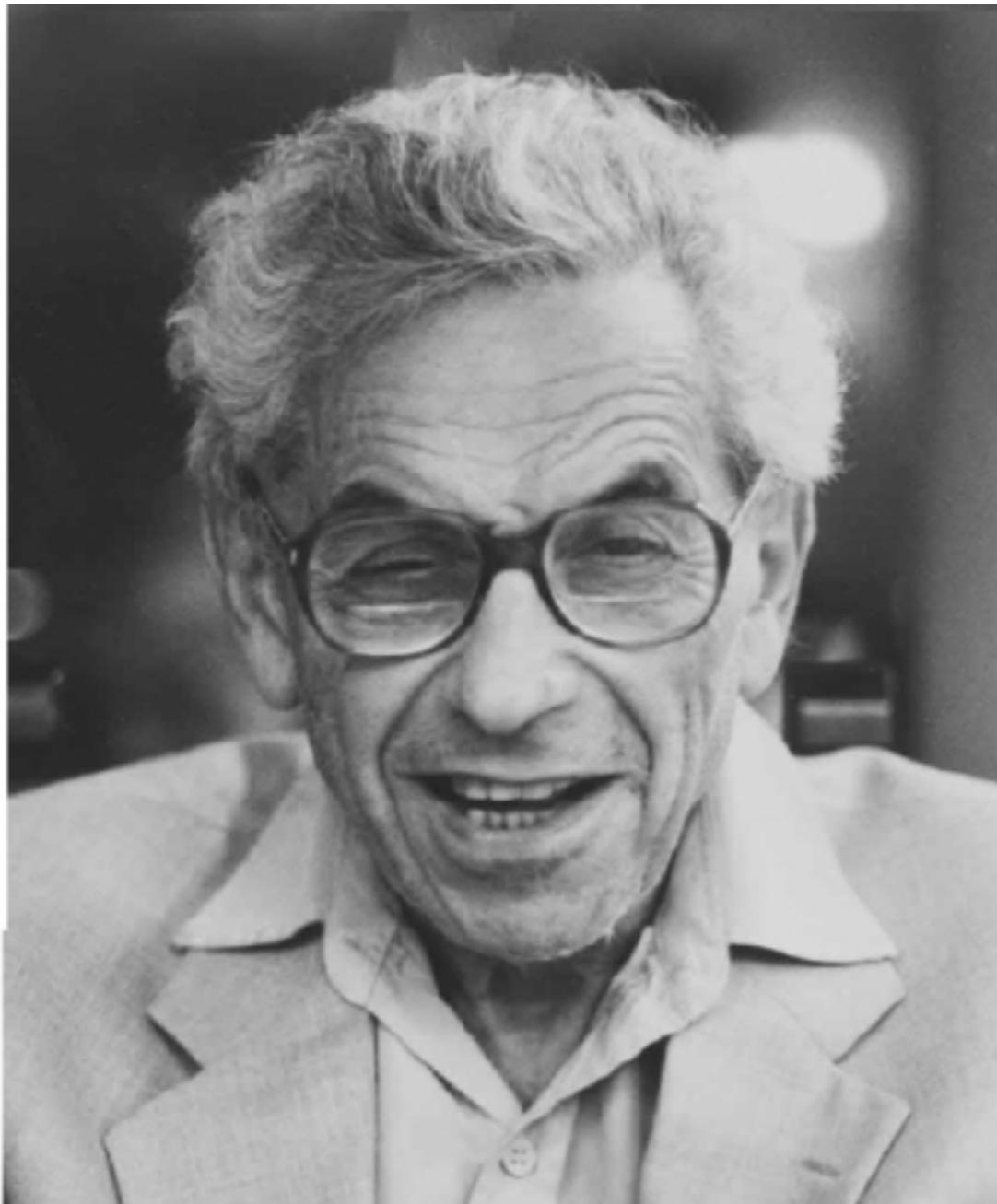
Proteins





Erdős-Renyi, 1959-1968

a.



Pál Erdős (1913-1996)

500+ co-authors

b.

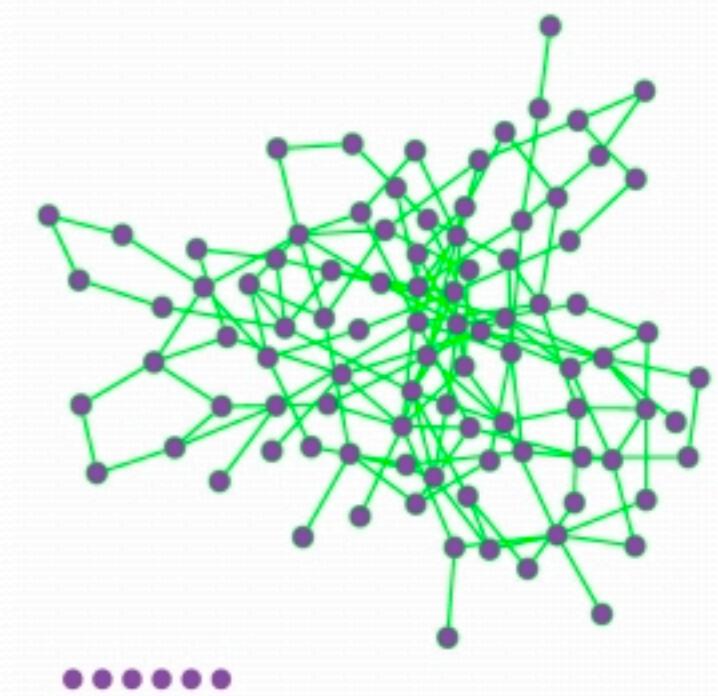
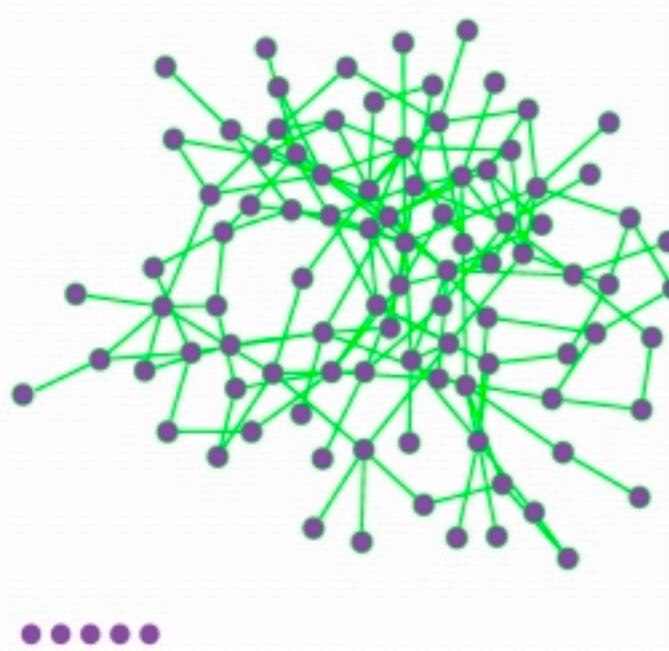
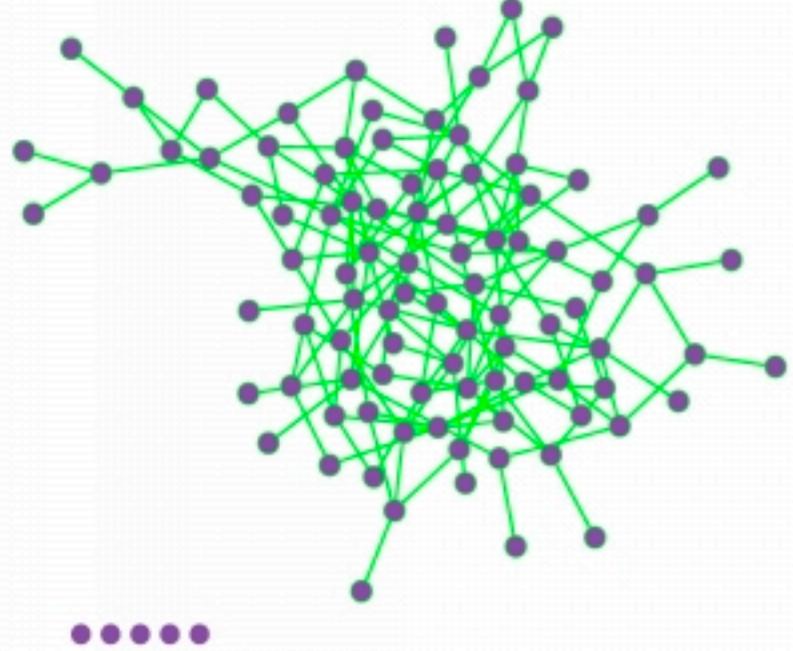
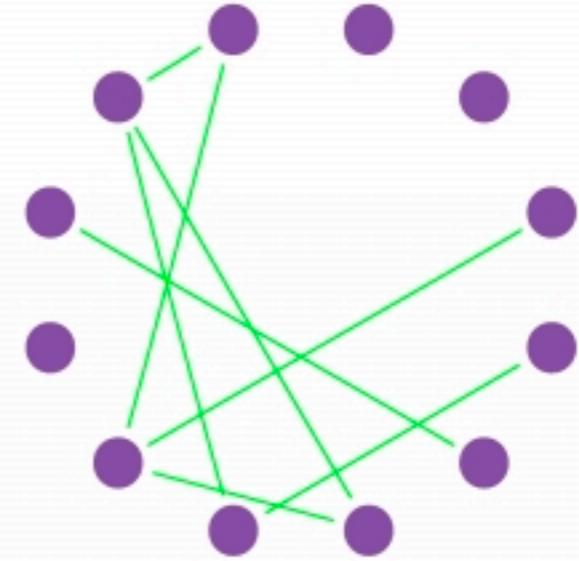
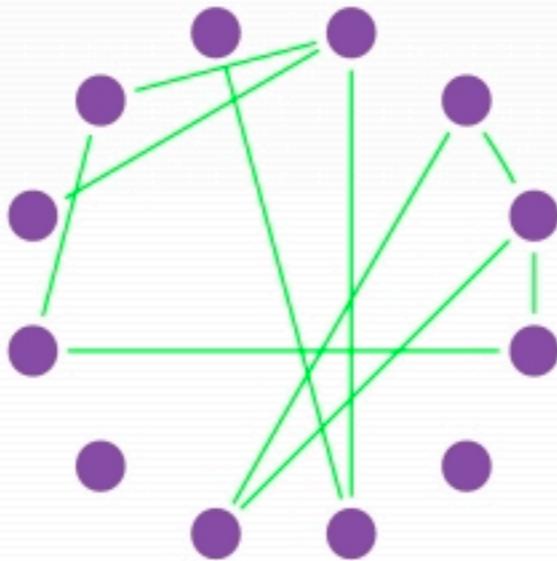
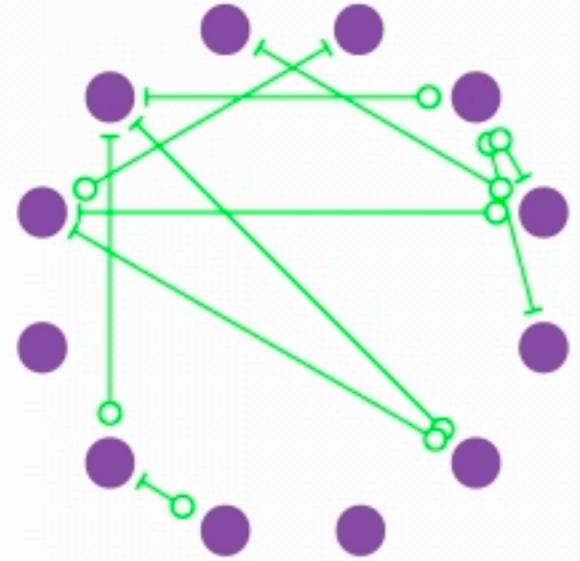


Alfréd Rényi (1921-1970)

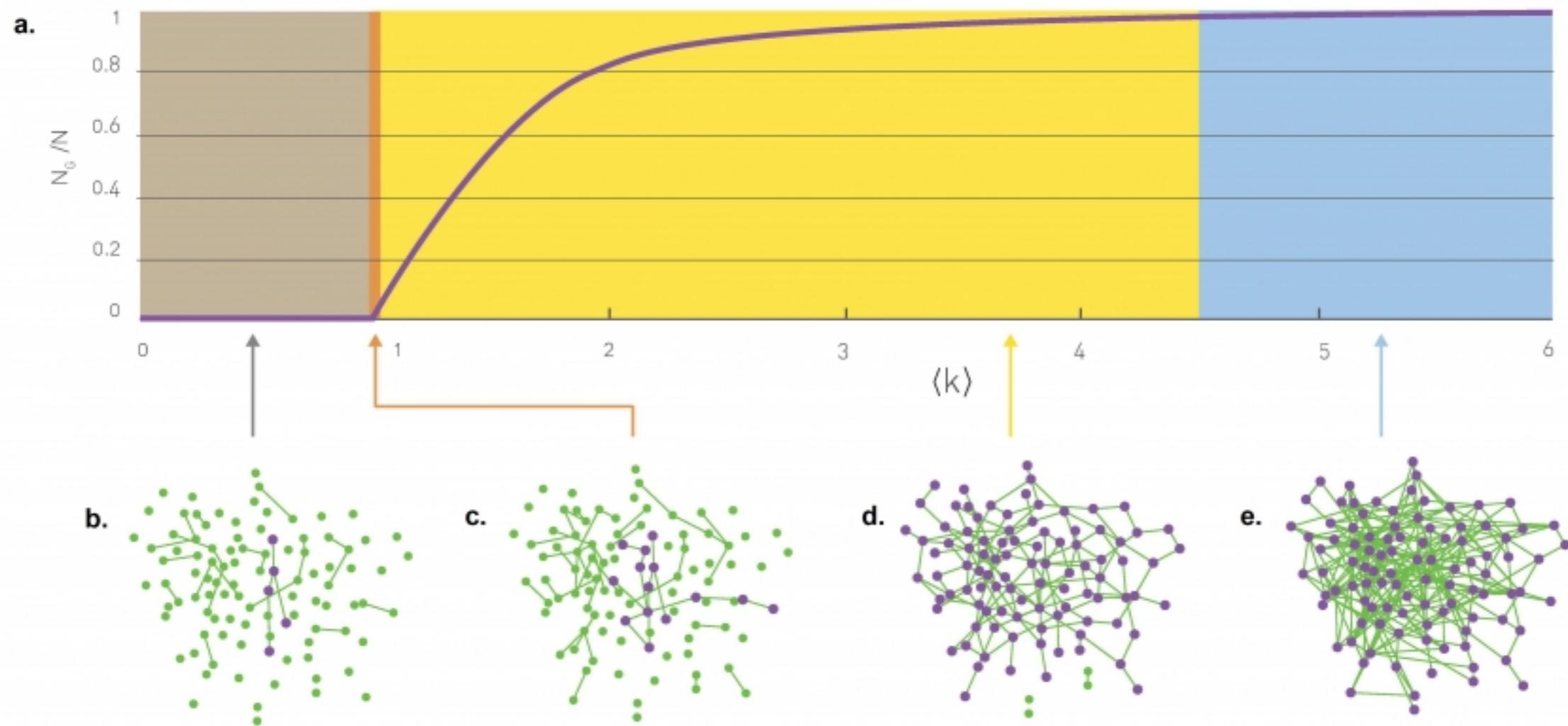
“A mathematician is a device for turning coffee into theorems”

ER RANDOM PROCESS

1. Create empty network with N isolated nodes
2. Select a node pair and generate a random number between 0 and 1. If the number exceeds p , connect the selected node pair with a link, otherwise leave them disconnected.
3. Repeat step (2) for each of the $N(N-1)/2$ node pairs.



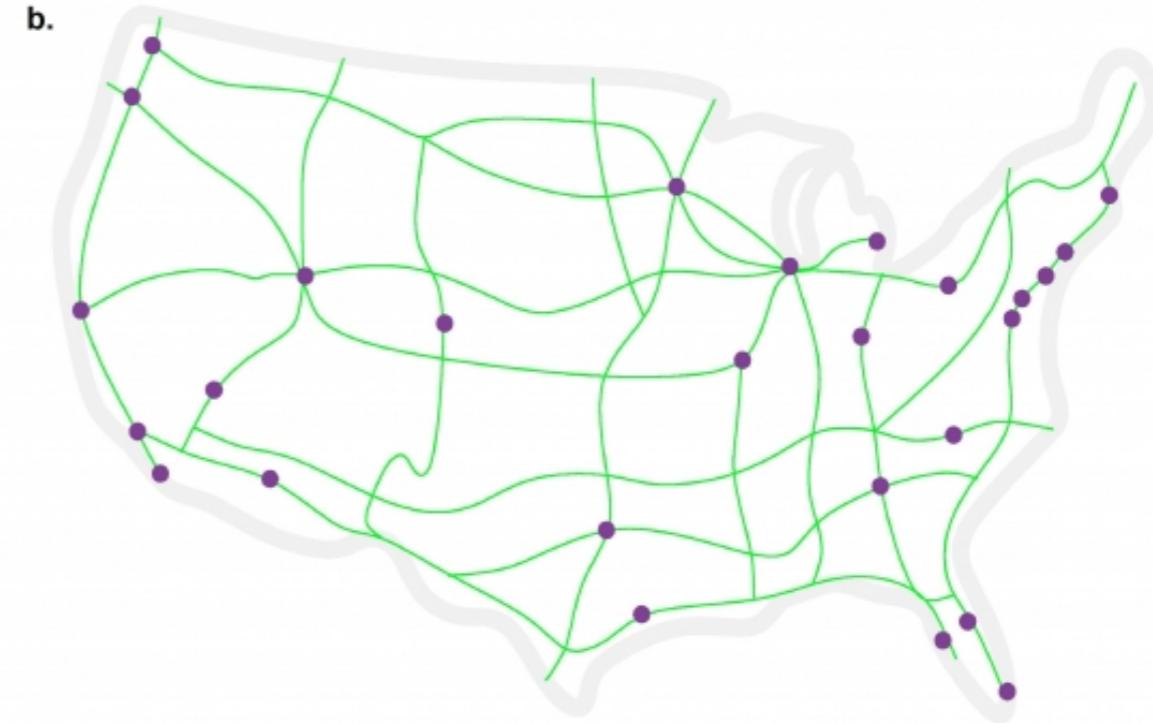
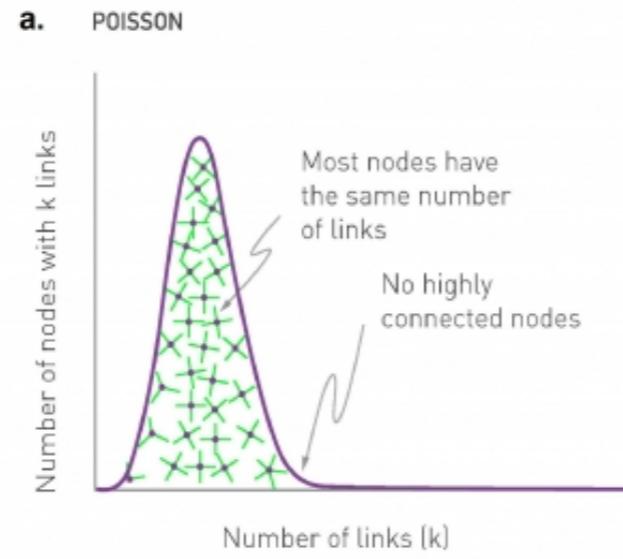
PERCOLATION TRANSITION



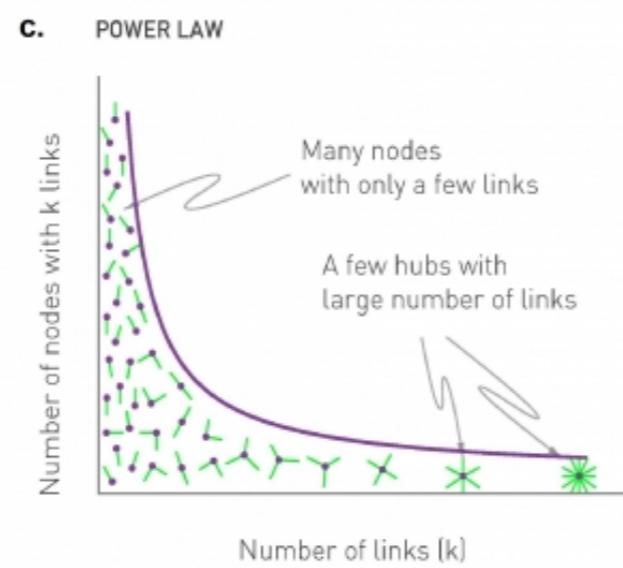
emergence of a giant component for growing density:
percolation transition at average degree $\langle k \rangle = 1$

REMINDER: SCALE-FREE NETWORKS

Erdos-Renyi
network

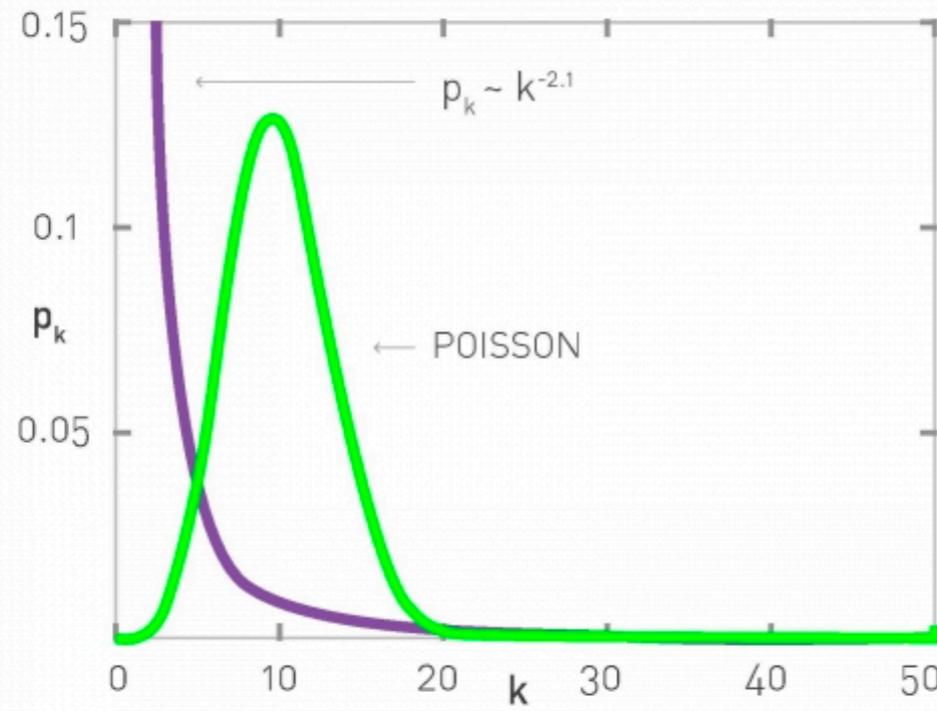


scale-free
network



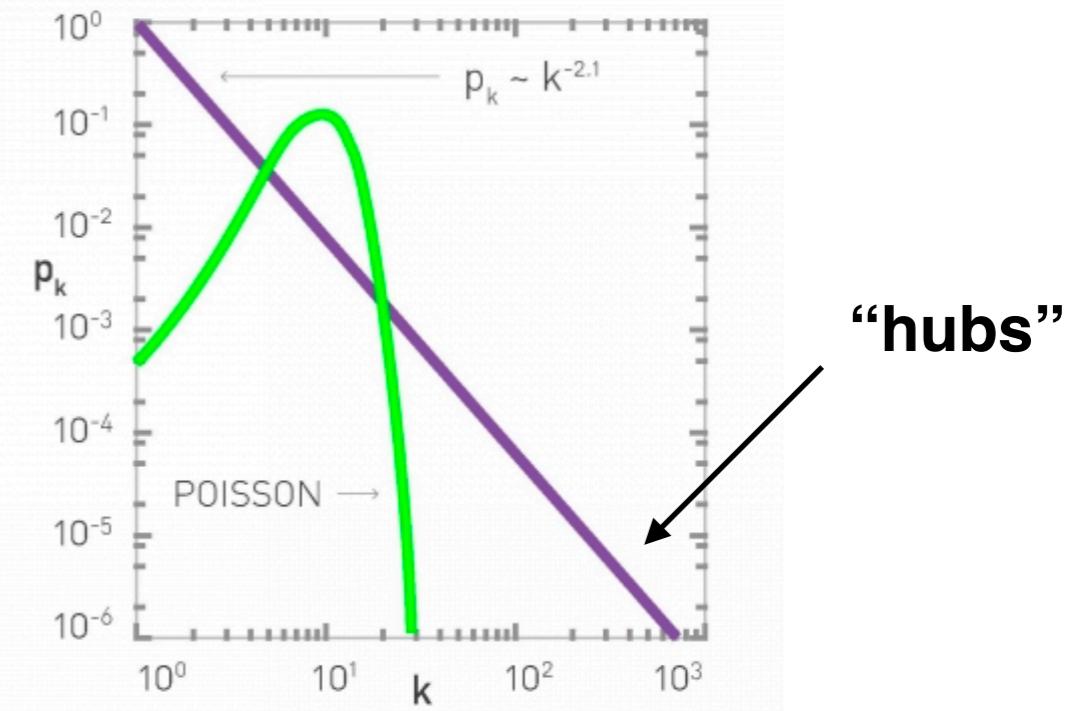
HUBS

a.



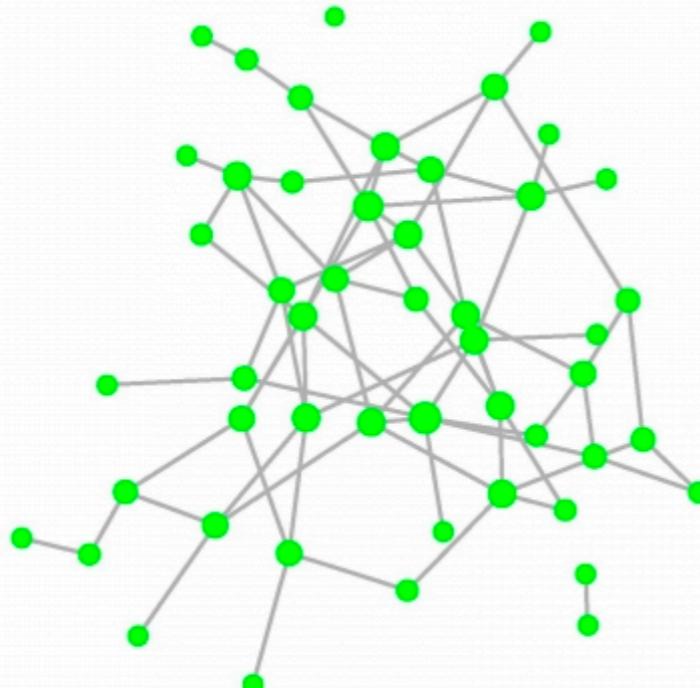
b.

log-log

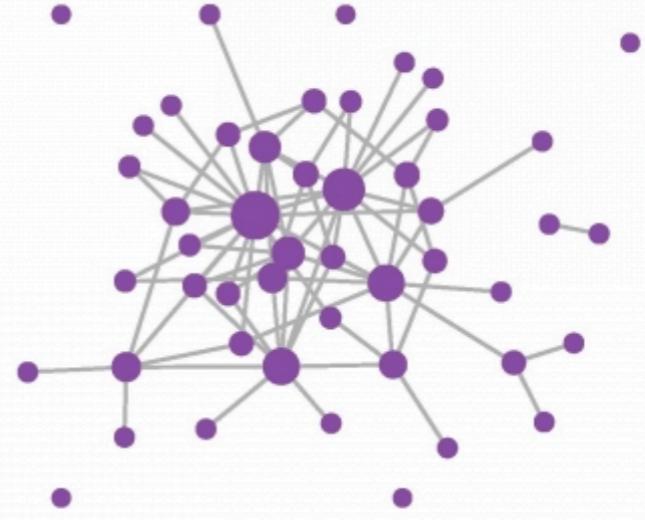


c.

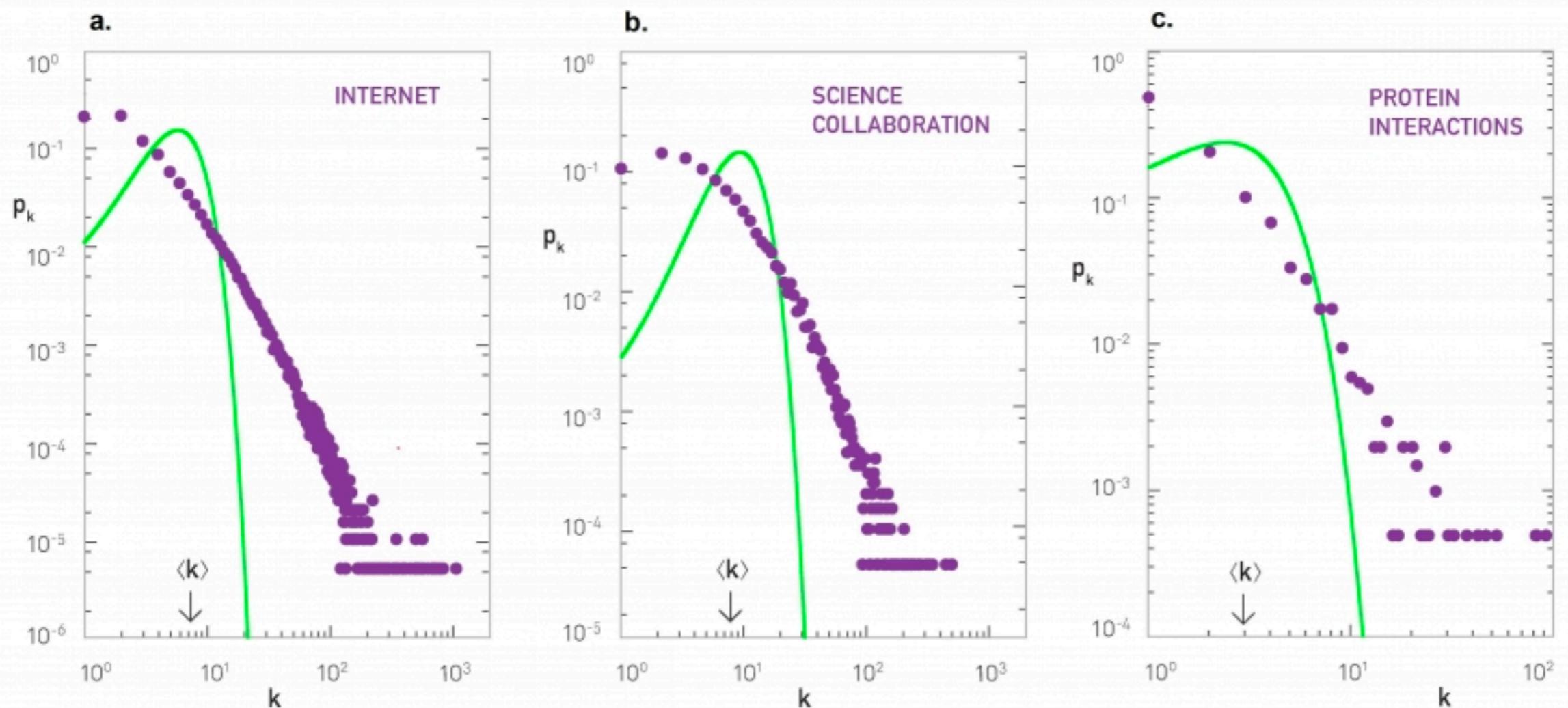
$$p_k \sim e^{-\langle k \rangle} \frac{\langle k \rangle^k}{k!}$$



d.

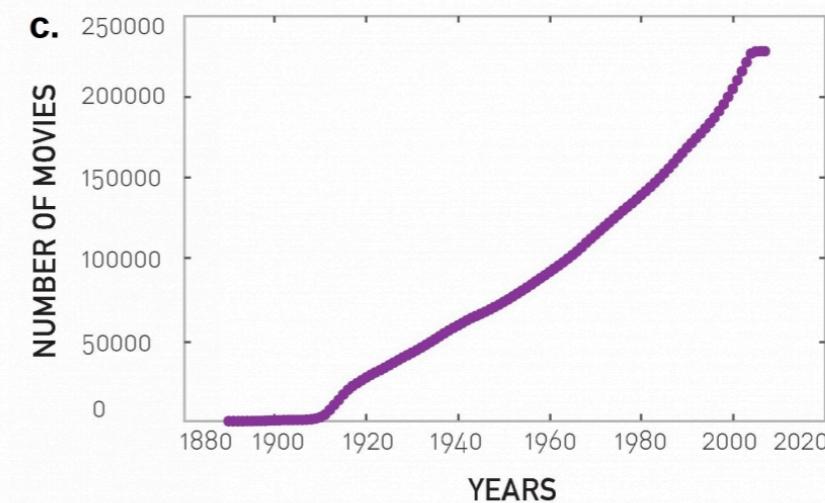
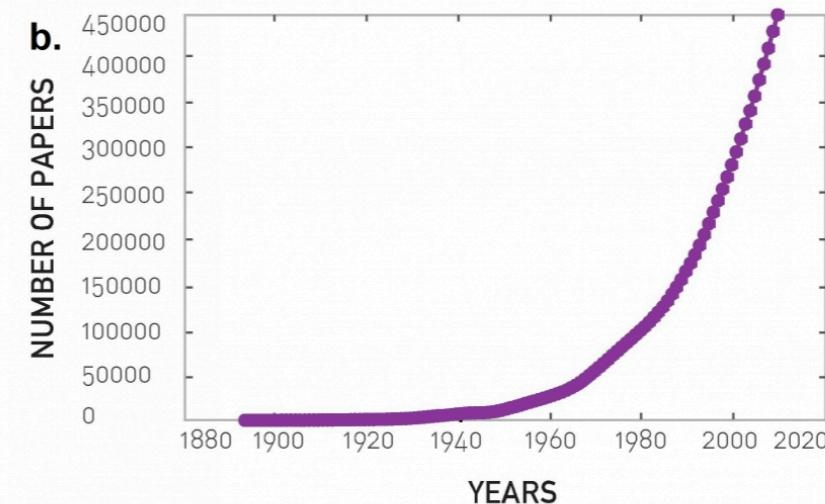
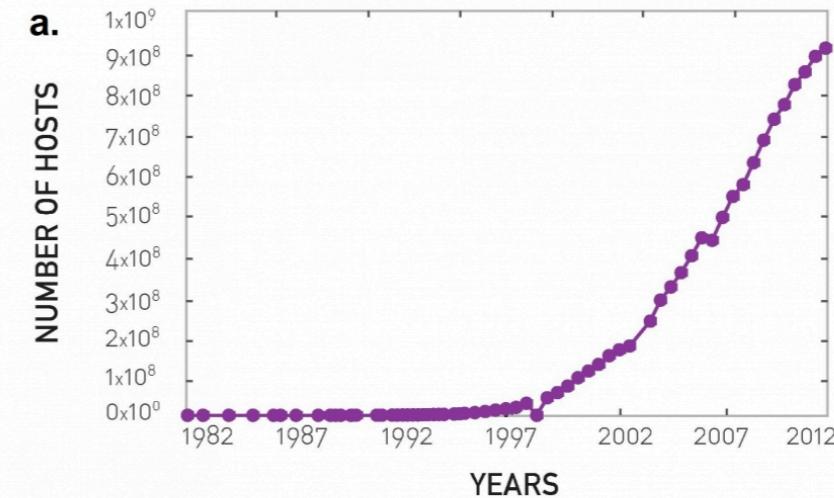


$$p(k) \sim k^{-\gamma}$$



REAL NETWORKS GROW IN TIME

number of nodes increases not
just density

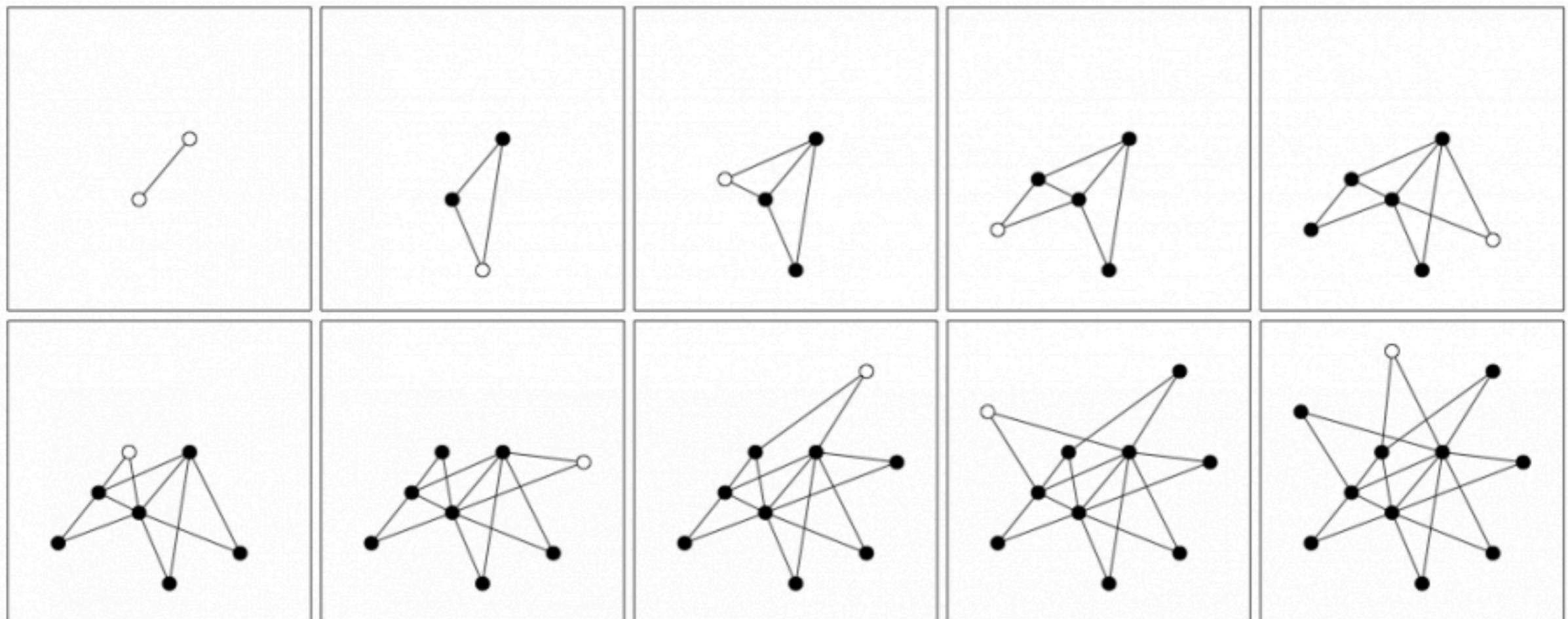


Barabasi-Albert model 1999

Preferential attachment

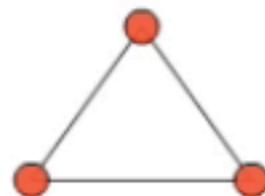
Probability $\Pi(k)$ that a link from a new node connects to node i :

$$\Pi(k_i) \sim \frac{k_i}{\sum_j k_k}$$



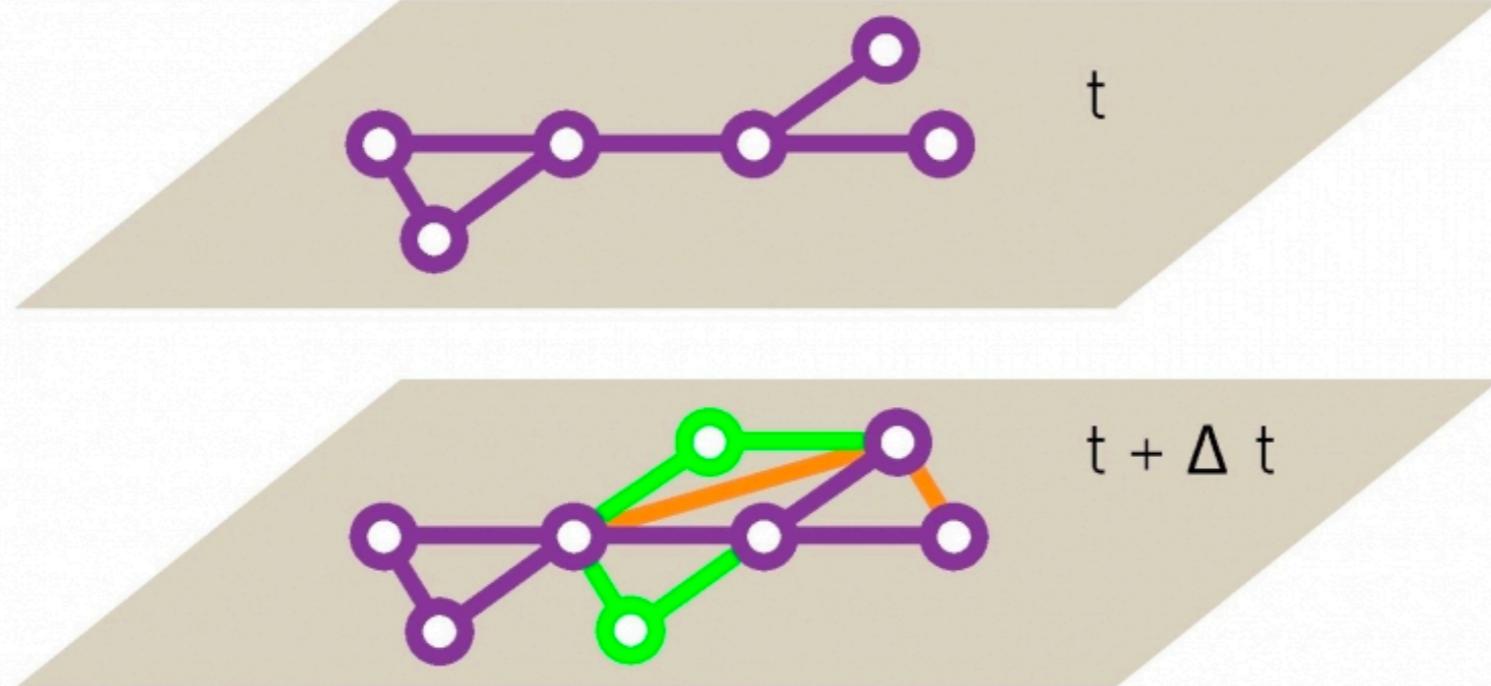
Barabasi-Albert model 1999

Preferential attachment

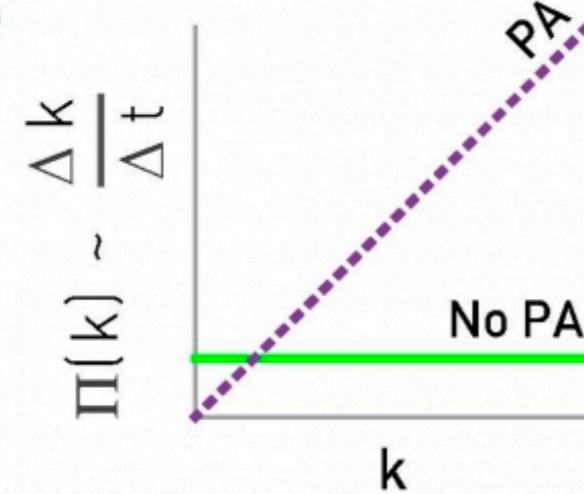


MEASURING PREFERENTIAL ATTACHEMENT (PA)

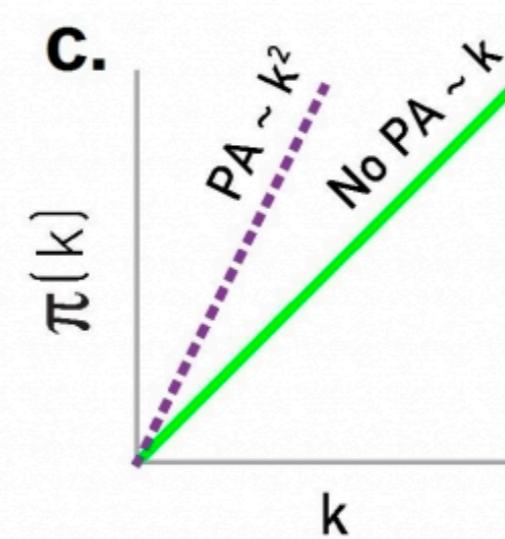
a.



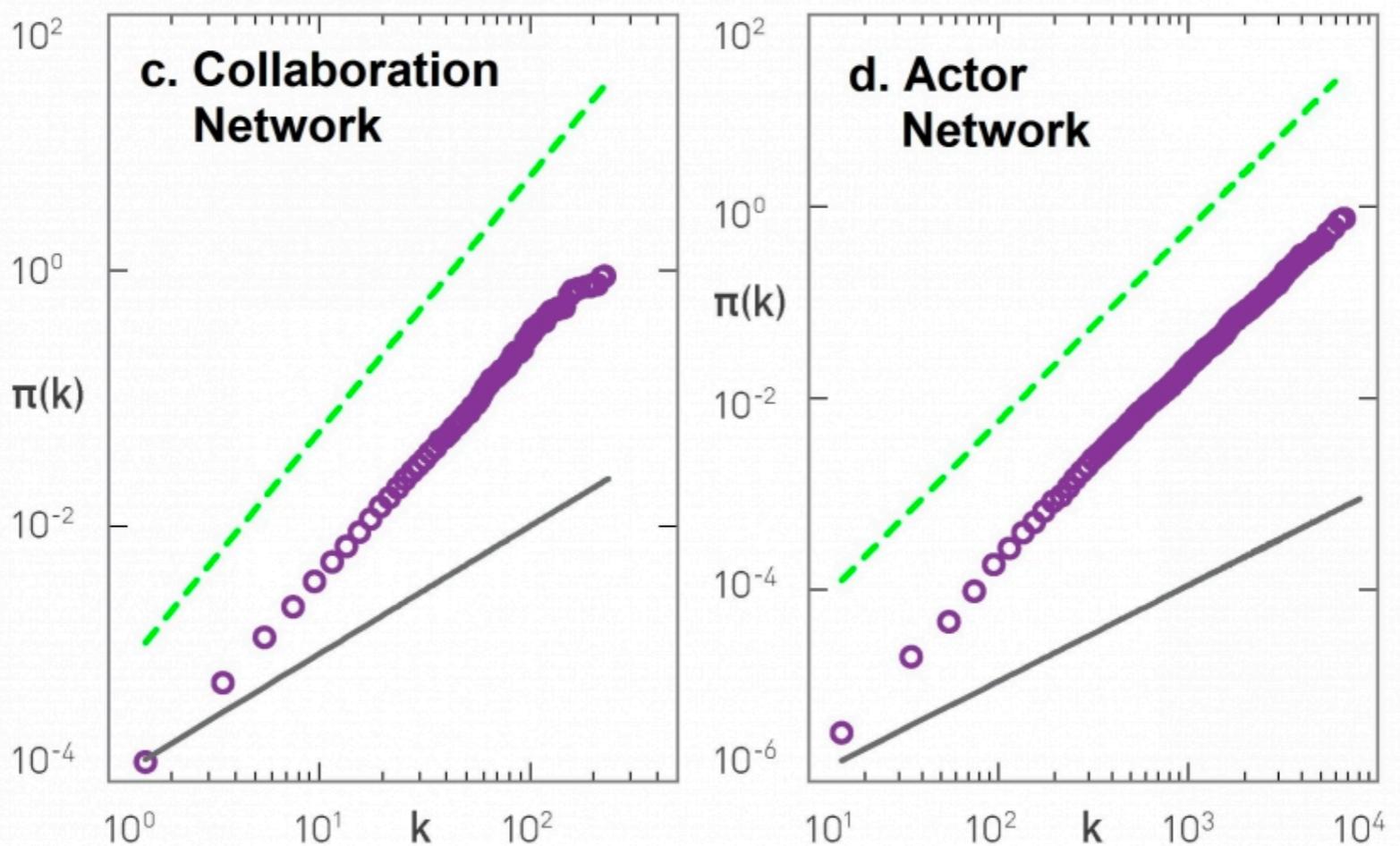
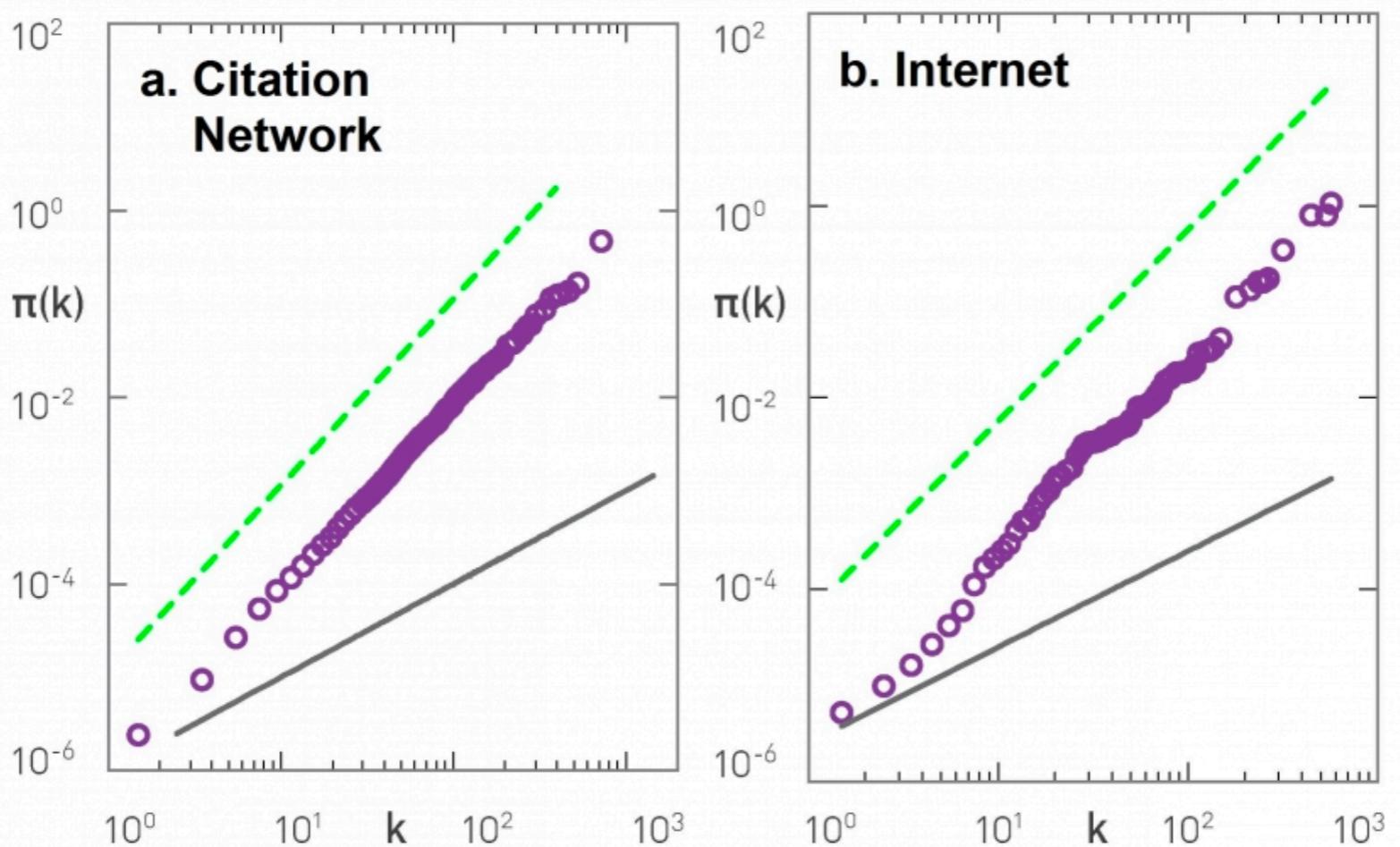
b.

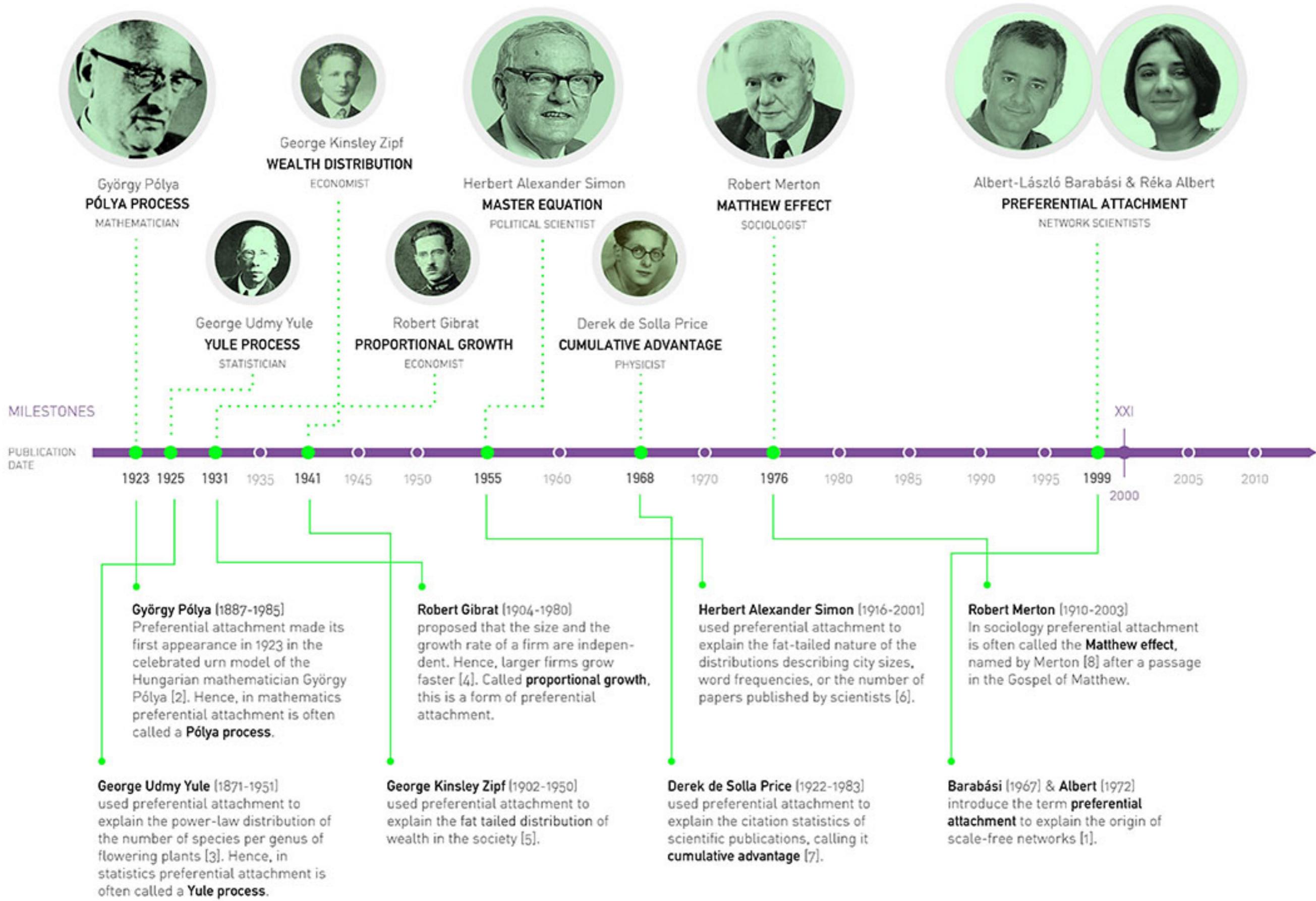


→



$$\pi(k) = \sum_{k_i=0}^k \Pi(k_i)$$





EXAMPLE IN BIOLOGY

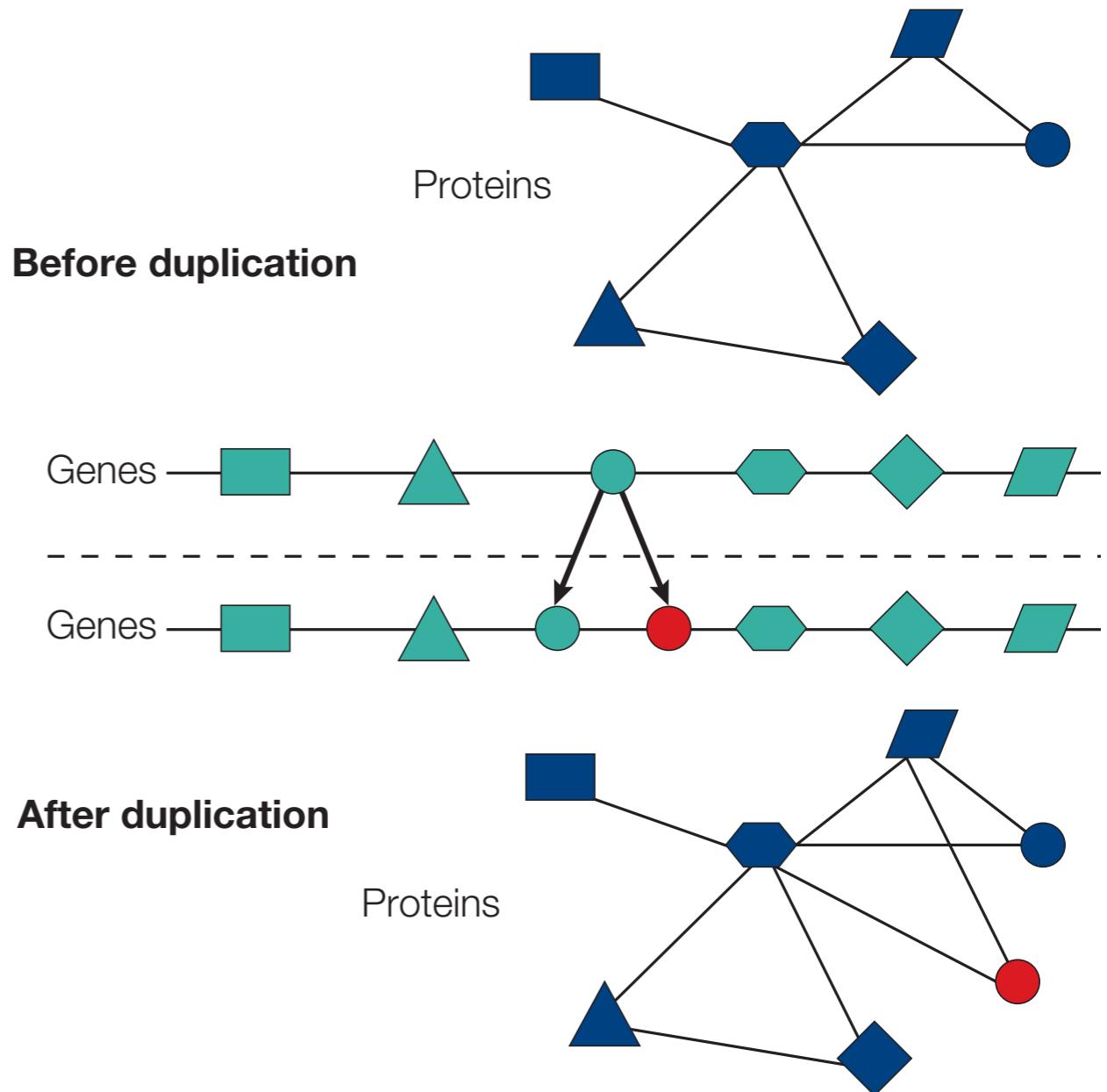
NETWORK BIOLOGY: UNDERSTANDING THE CELL'S FUNCTIONAL ORGANIZATION

Albert-László Barabási* & Zoltán N. Oltvai†

A key aim of postgenomic biomedical research is to systematically catalogue all molecules and their interactions within a living cell. There is a clear need to understand how these molecules and the interactions between them determine the function of this enormously complex machinery, both in isolation and when surrounded by other cells. Rapid advances in network biology indicate that cellular networks are governed by universal laws and offer a new conceptual framework that could potentially revolutionize our view of biology and disease pathologies in the twenty-first century.

Barabasi et al., Nat Rev Genet 2004

preferential attachment: highly connected proteins have higher chance to be connected to a duplicated protein



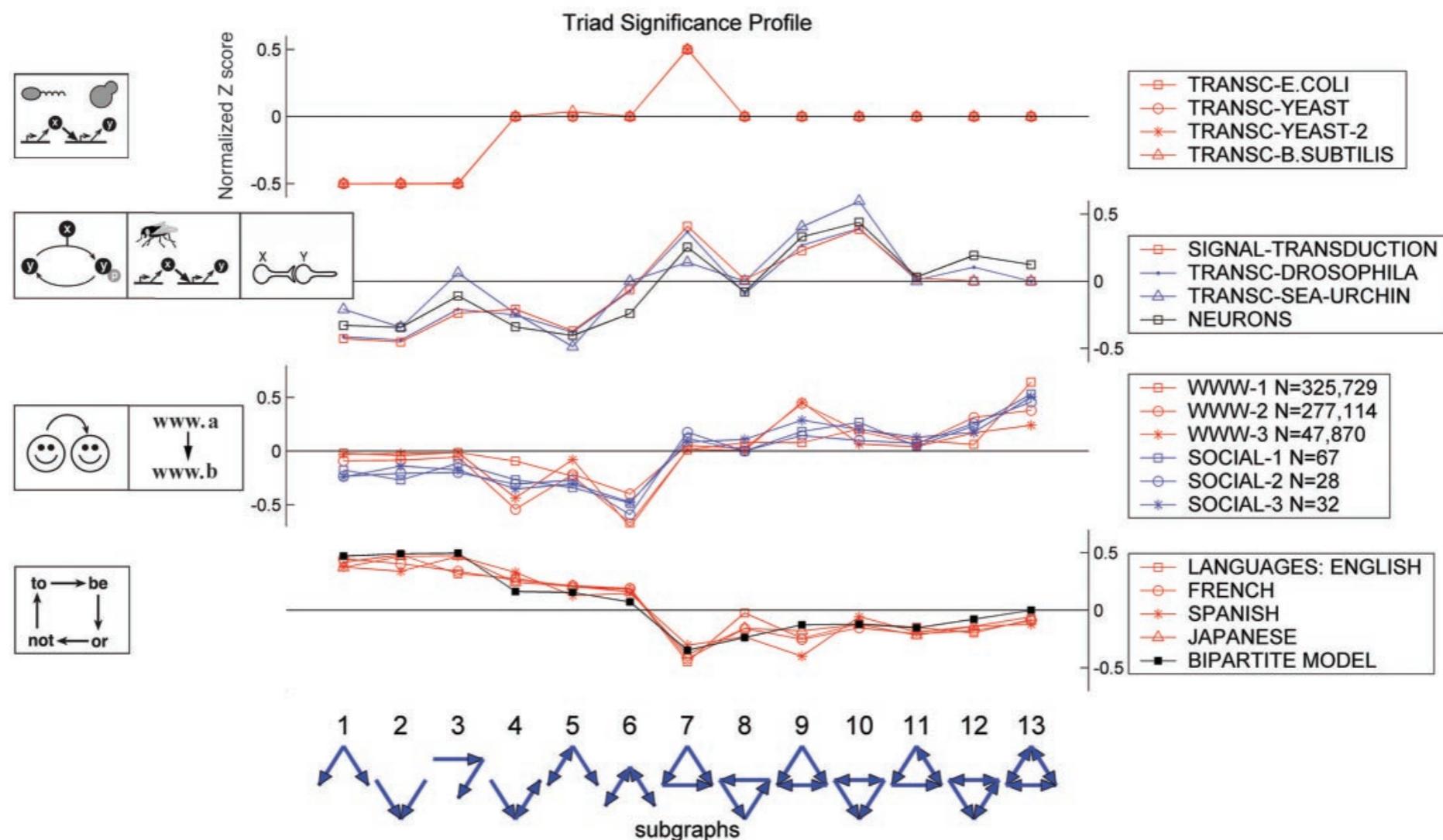
HIGHER ORDERS: TRIADIC CLOSURE

Superfamilies of Evolved and Designed Networks

Ron Milo, Shalev Itzkovitz, Nadav Kashtan, Reuven Levitt,
Shai Shen-Orr, Inbal Ayzenstiel, Michal Sheffer, Uri Alon*

5 MARCH 2004 VOL 303 SCIENCE

in real networks, there is also a preferential attachment to form certain **network motifs** (like the triadic “feed-forward loop”)



ADVANCED METHODS

Exponential Random Graph Models (ERGMs)

(eg in R https://statnet.github.io/Workshops/ergm_tutorial.html)

$$p(X = x) = \frac{\exp\{\theta'z(x)\}}{\kappa(\theta)}$$

Where:

X is a random network on n nodes

x is the observed network

θ is a vector of parameters (like regression coefficients)

$z(x)$ is a vector of network statistics

κ is a normalizing constant, to ensure the probabilities sum to 1:

$$\kappa(\theta) = \sum_{\substack{\text{all possible} \\ \text{graphs} \\ x}} \exp\{\theta'z(x)\}$$

Typical terms in the model

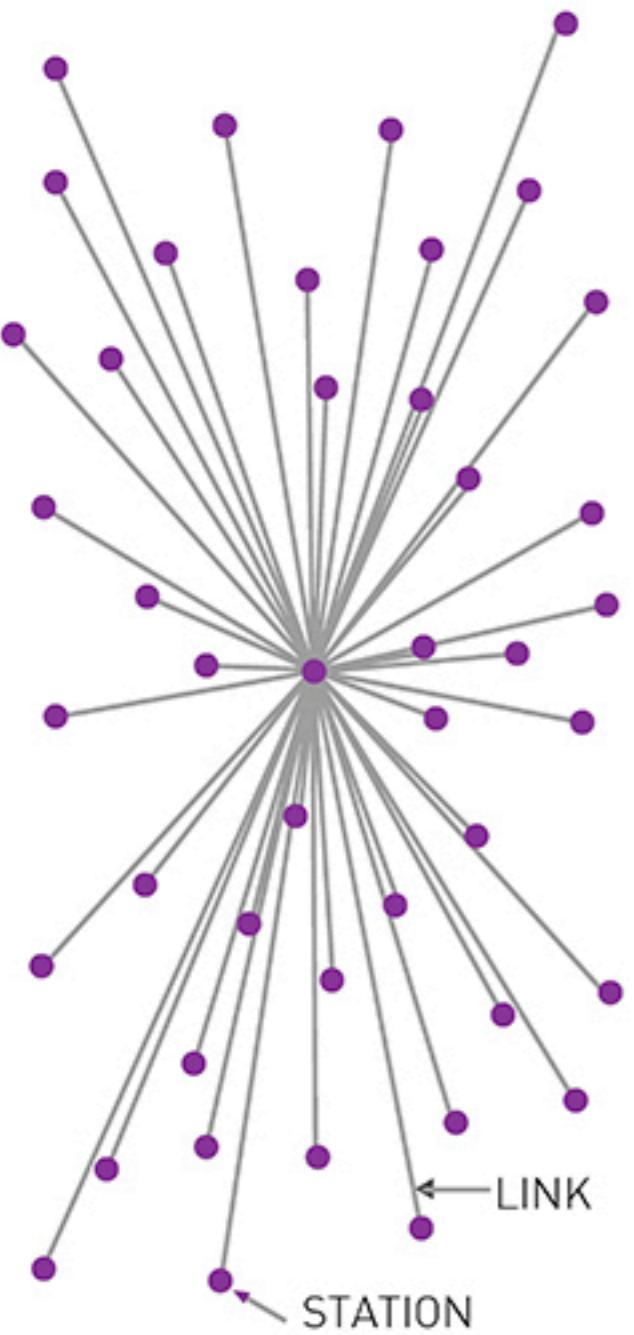
- Edges
- Mutuality
- Homophily
- Degree distribution
- Triangles
- ...

HOW TO ATTACK A NETWORK

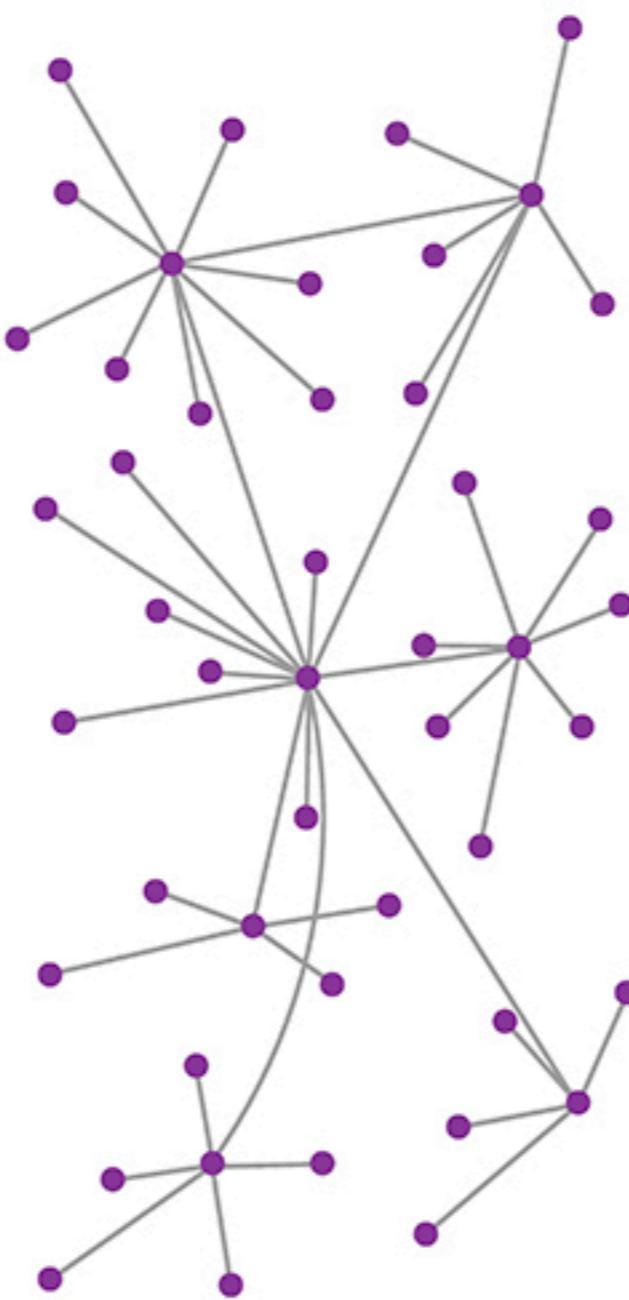
ROBUSTNESS AND FRAGILITY

C

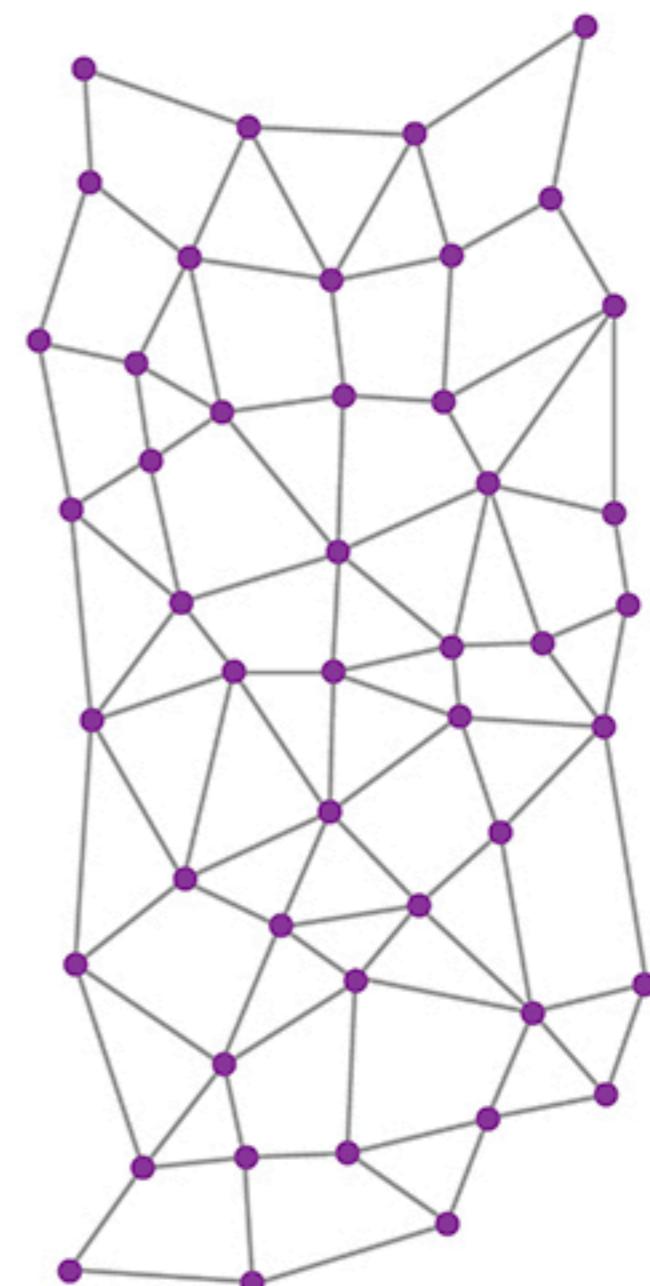




a. CENTRALIZED



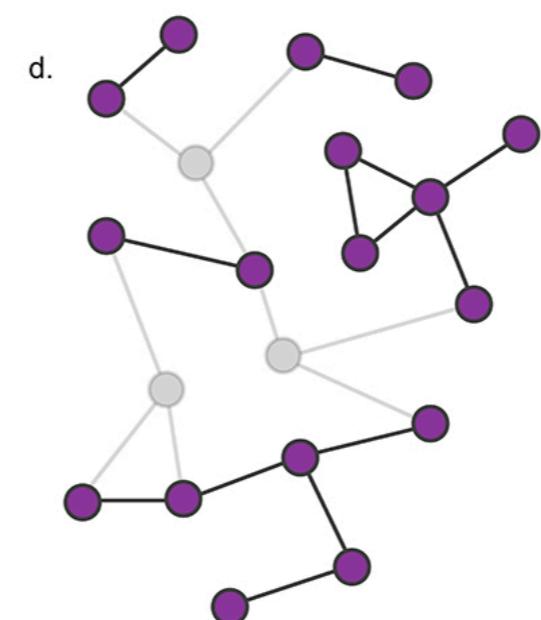
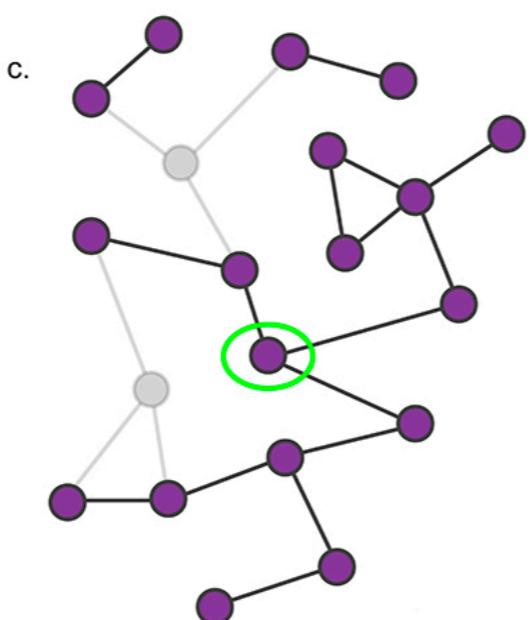
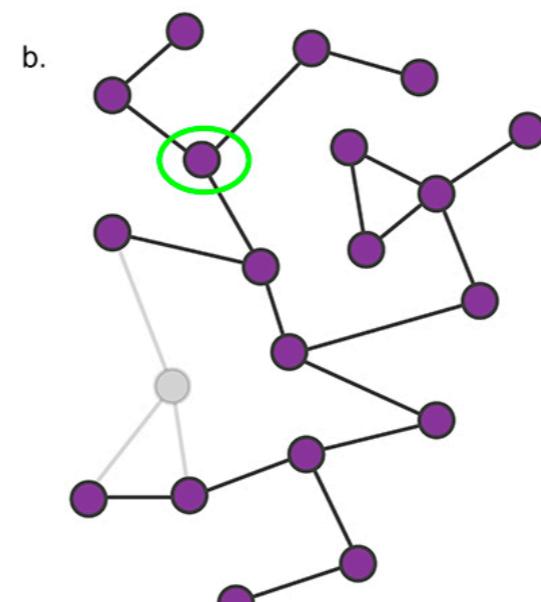
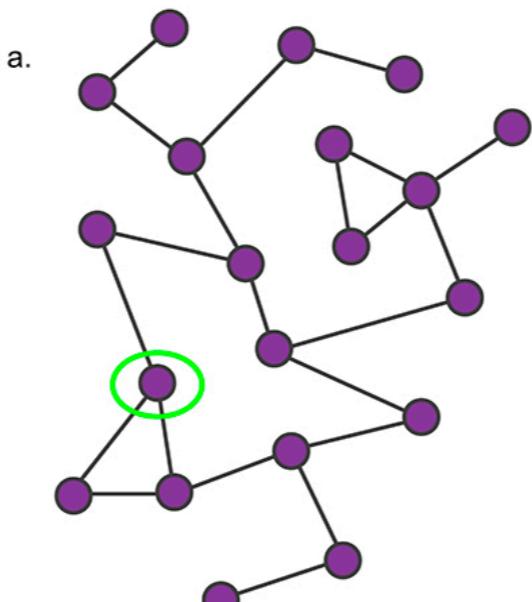
b. DECENTRALIZED



c. DISTRIBUTED

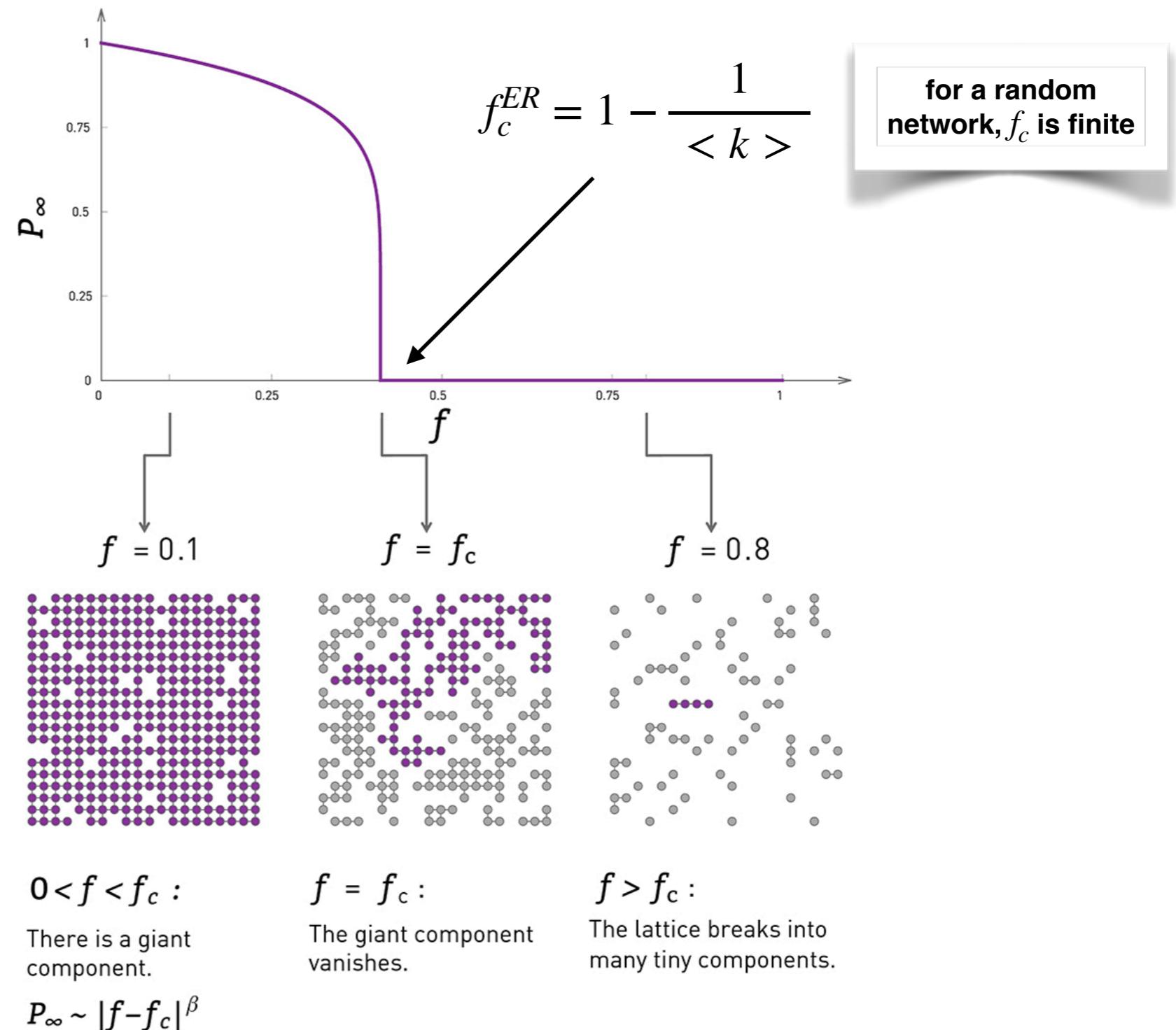
RANDOM FAILURES

At each step, remove a node at random along with its links



INVERSE PERCOLATION

Proportion of nodes in
the “giant” (largest
connected) component

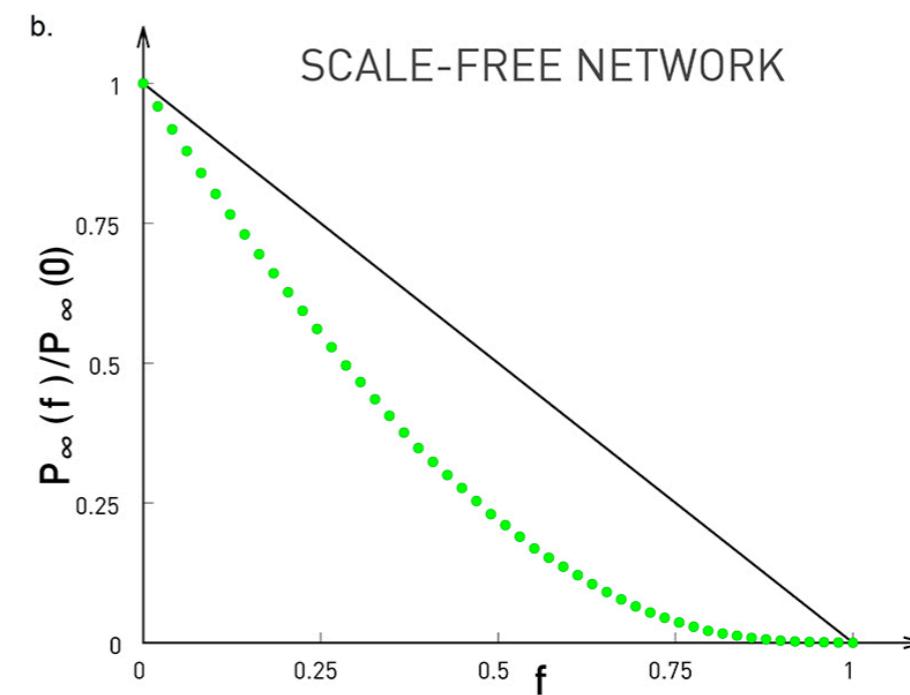
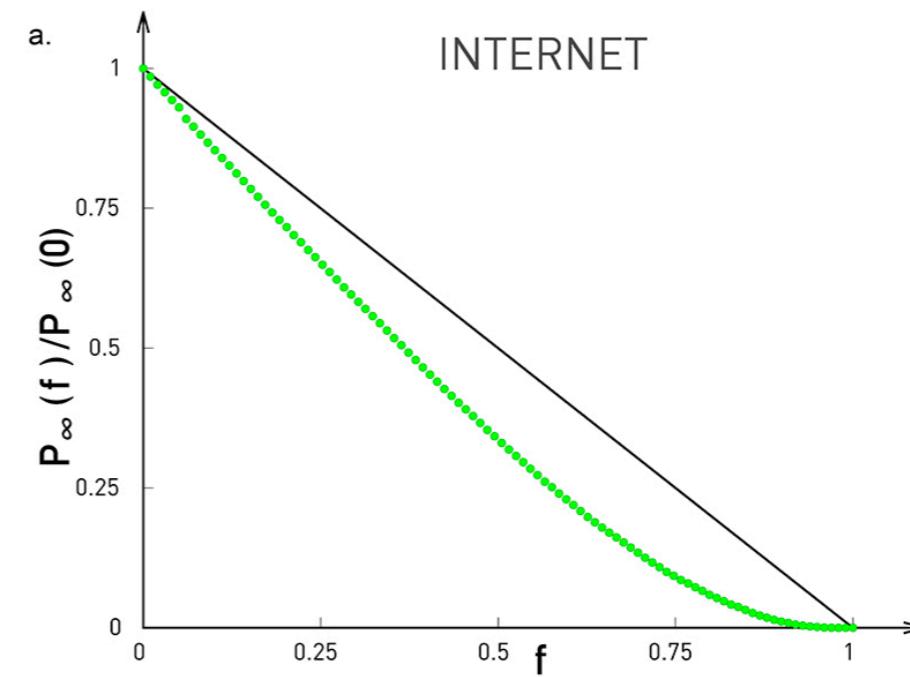


ROBUSTNESS OF SCALE FREE NETWORKS

general solution:

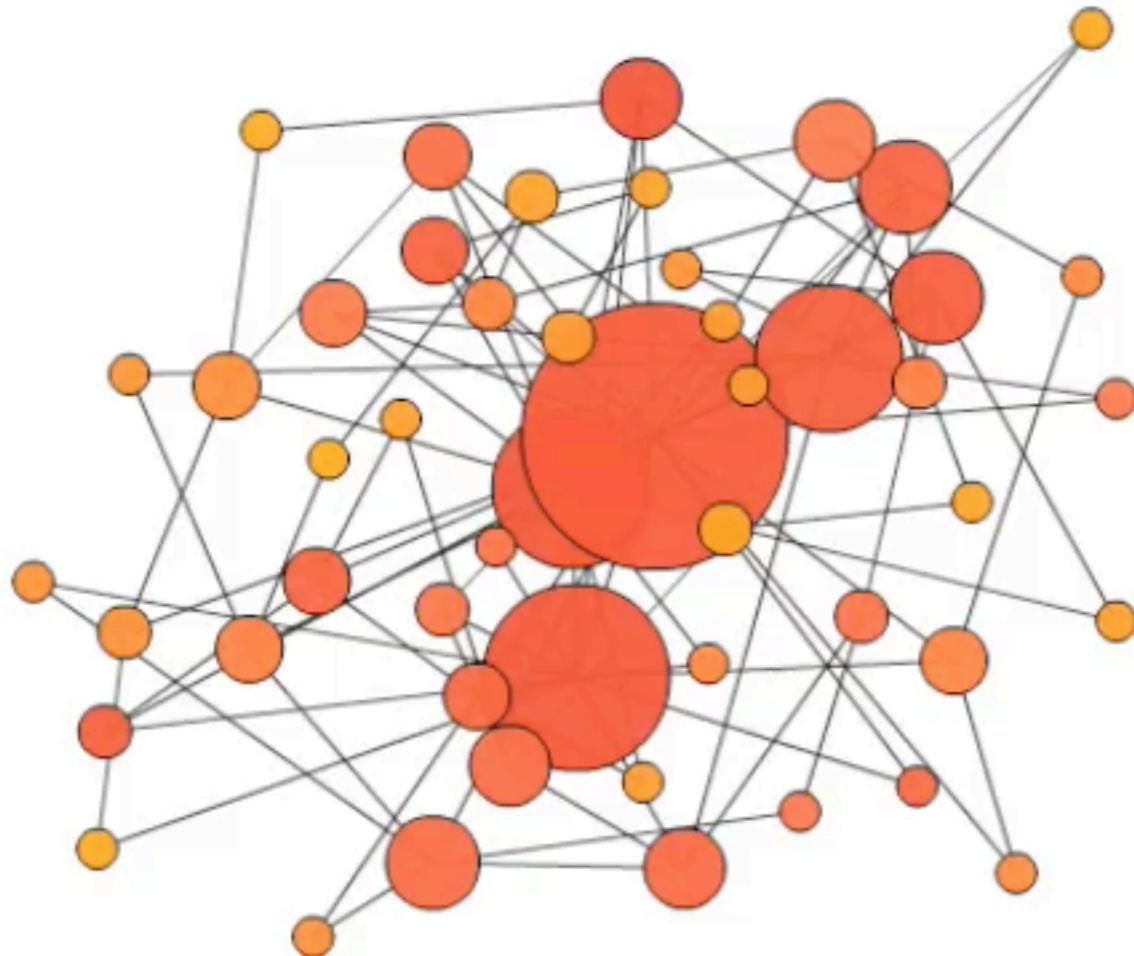
$$f_c = 1 - \frac{1}{\frac{\langle k^2 \rangle}{\langle k \rangle} - 1}$$

for many scale-free networks,
 $\langle k^2 \rangle \rightarrow \infty$ to that $f_c = 1$ and the
network is absolutely robust!



Scale-free Network Under Node Failures

random node removal



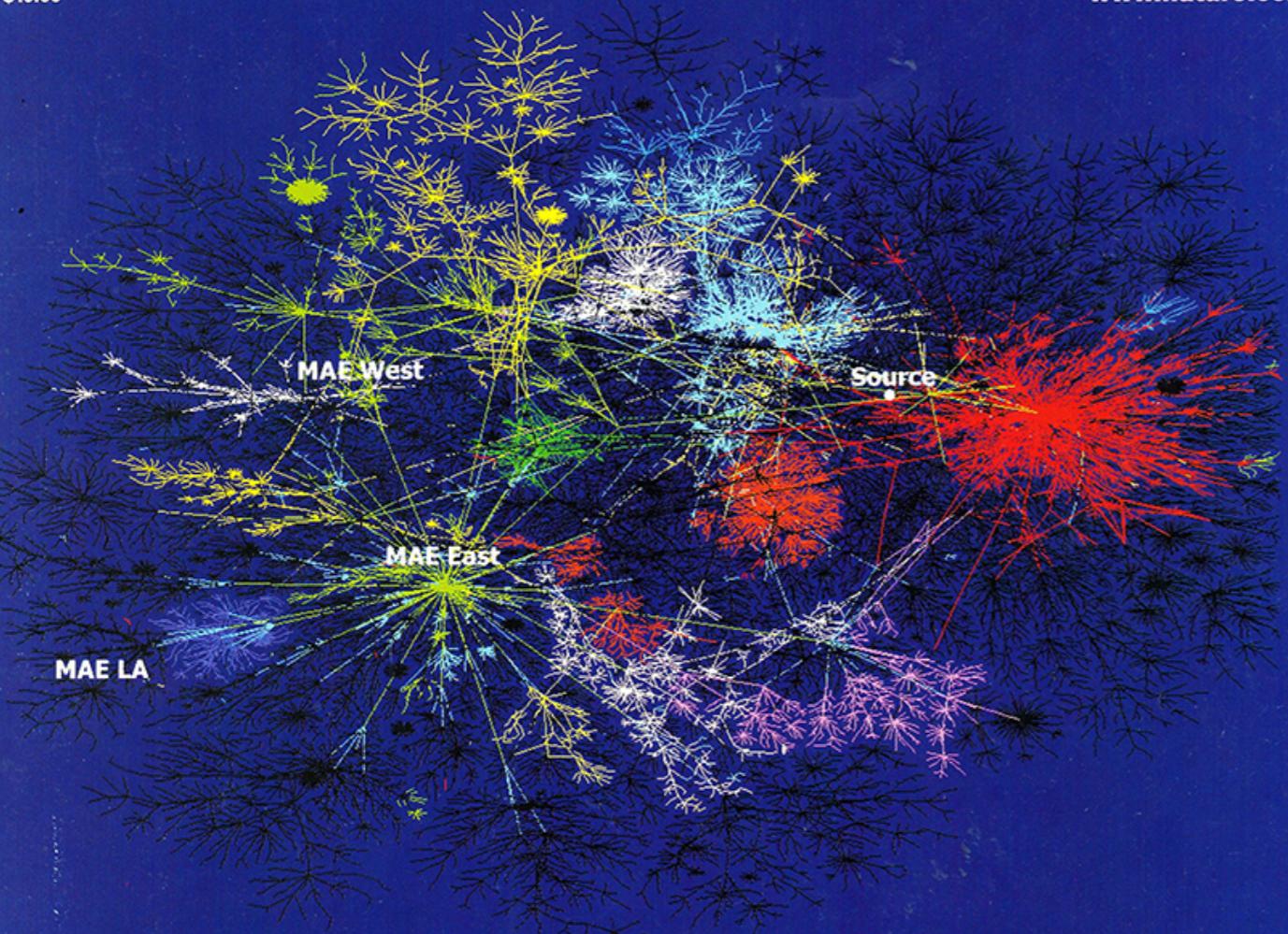
27 July 2000

International weekly journal of science

nature

\$10.00

www.nature.com



Achilles' heel of the Internet

Obesity Mice that eat more but weigh less

Ocean anoxic events Not all at sea

Cell signalling Fringe sweetens Notch

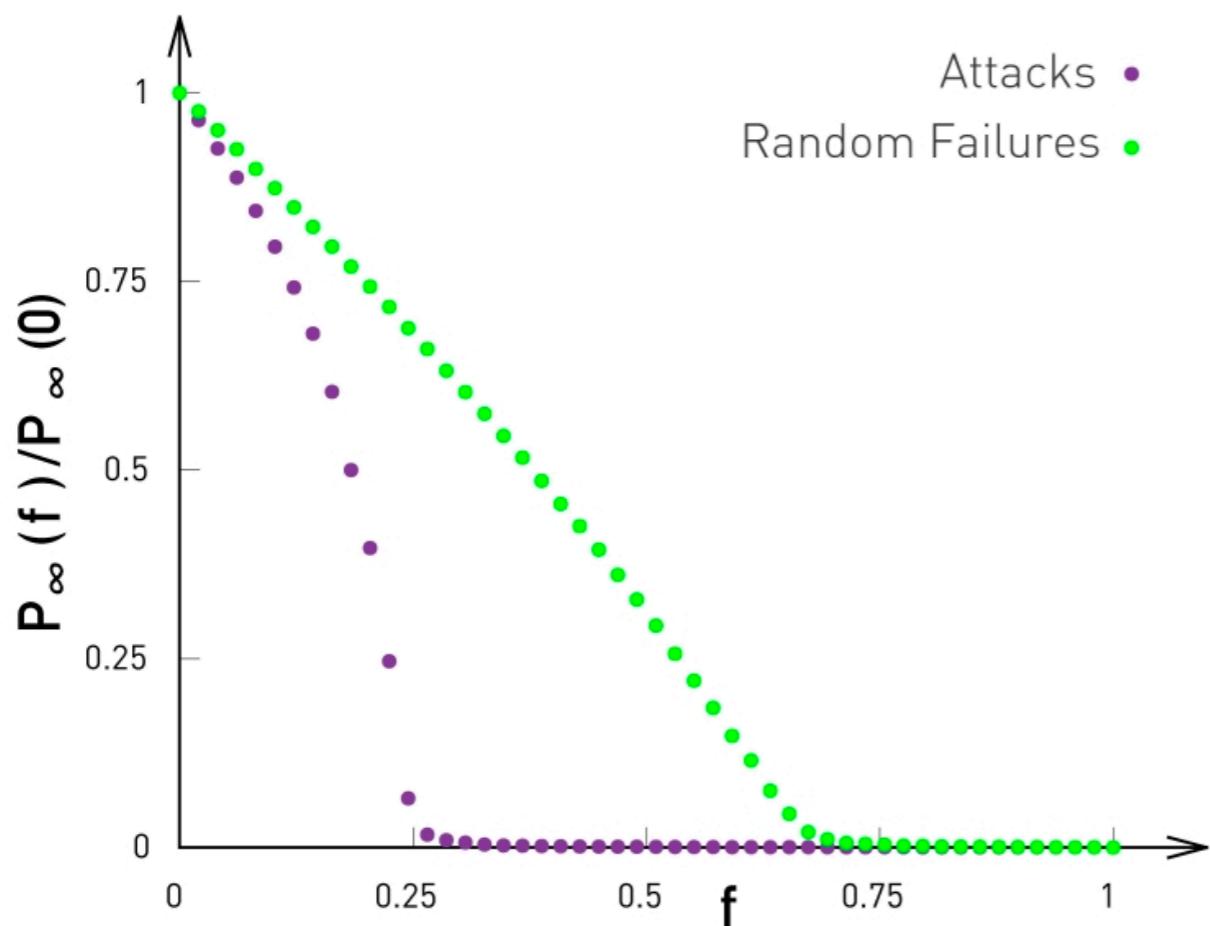
Albert,R., Jeong,H. and Barabási,A. (2000)
Error and attack tolerance of complex
networks. *Nature*, 406, 378–382.

new on the market

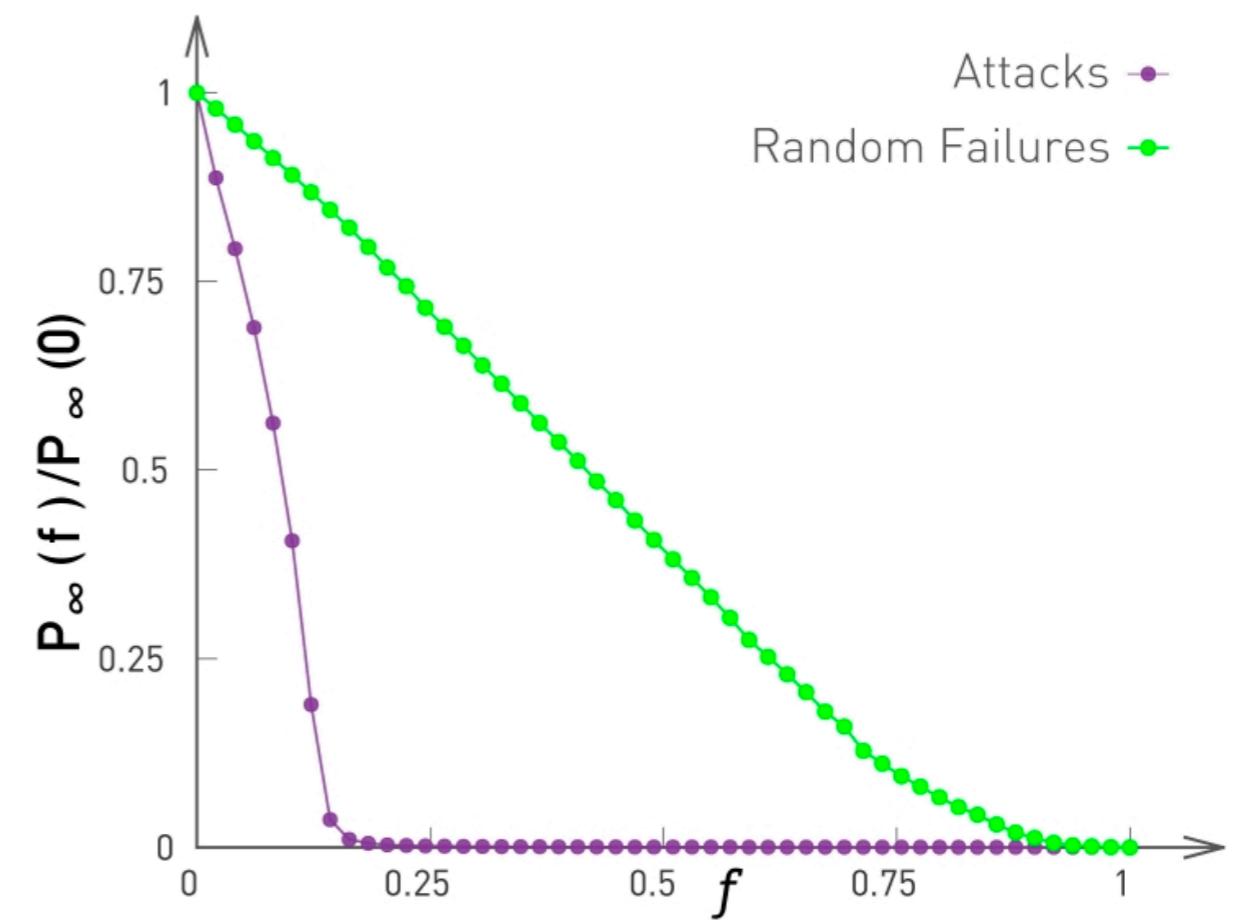
oligonucleotides

NETWORKS UNDER ATTACK

Random (ER) network

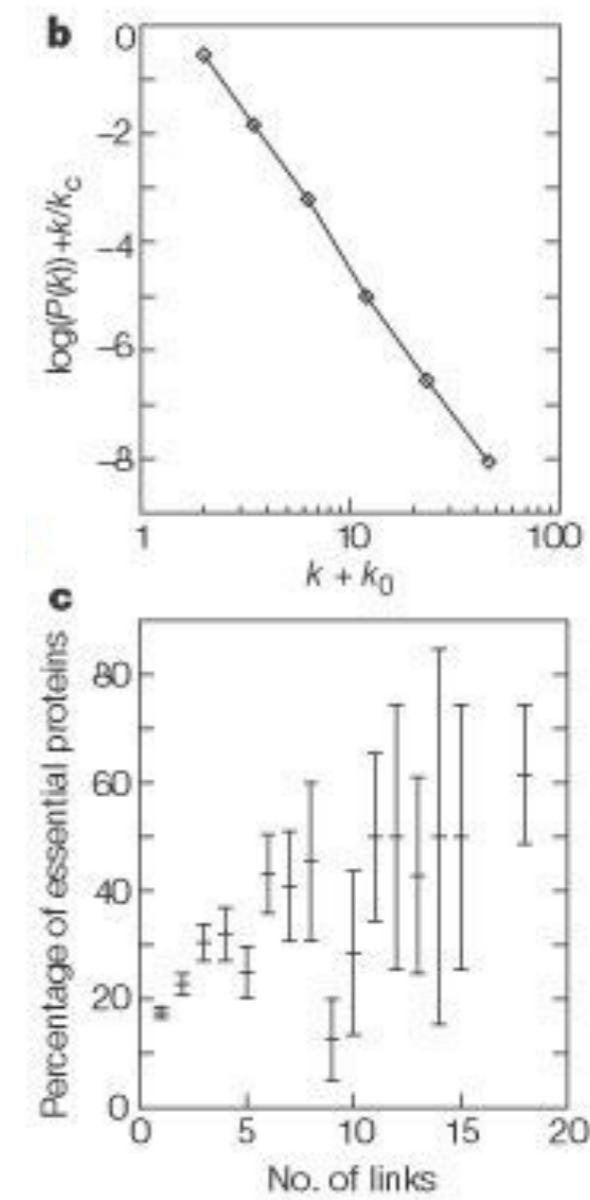
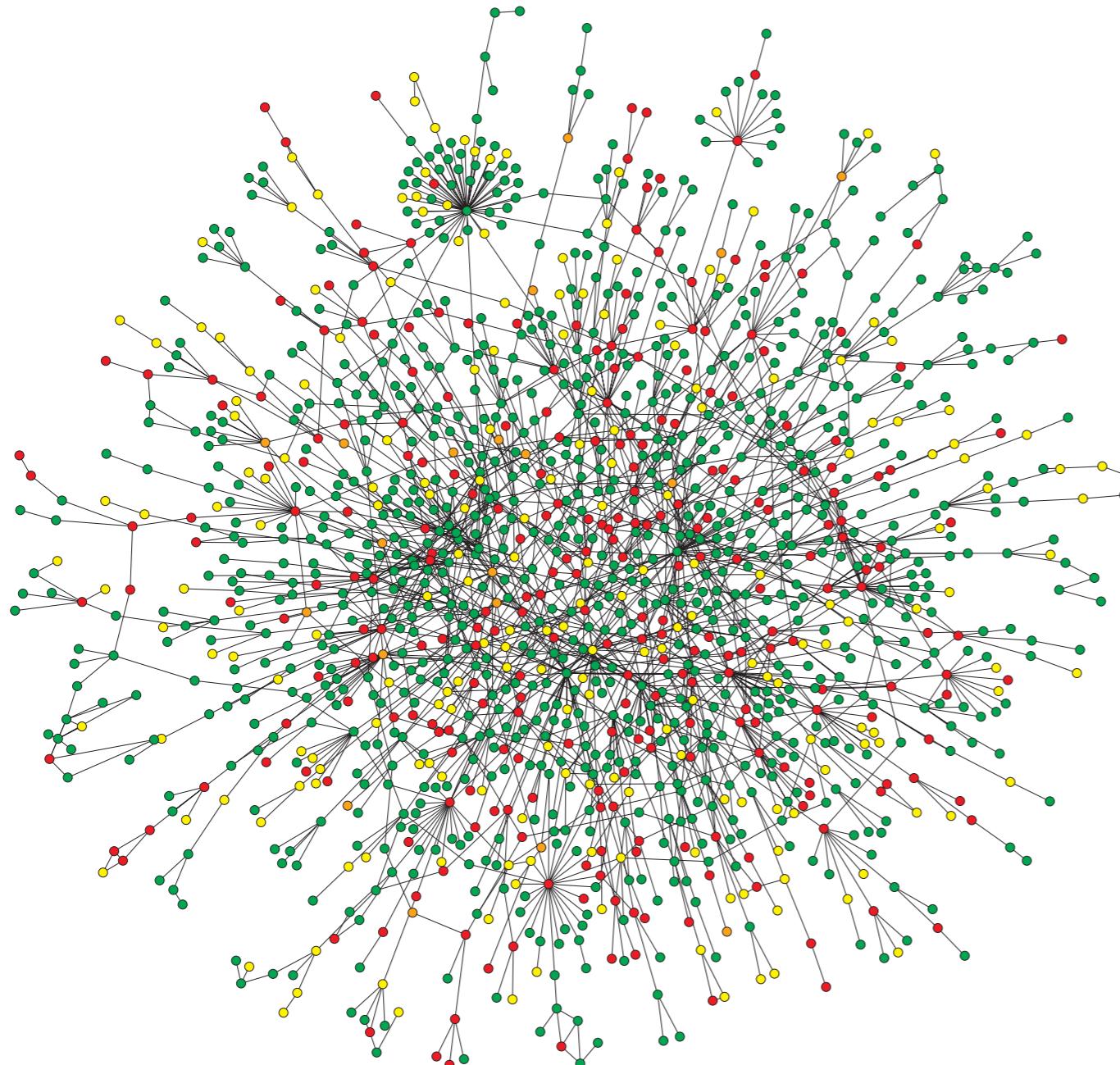


Scale-free (BA) network



while very robust against random attacks, scale-free networks
are very fragile against targeted attacks on hubs!

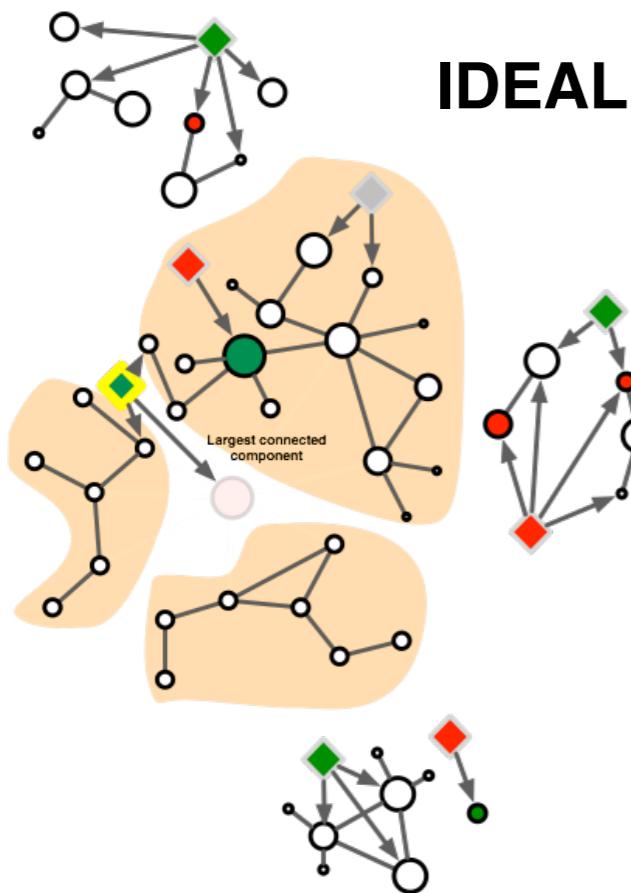
EXAMPLE: CENTRALITY-LETHALITY



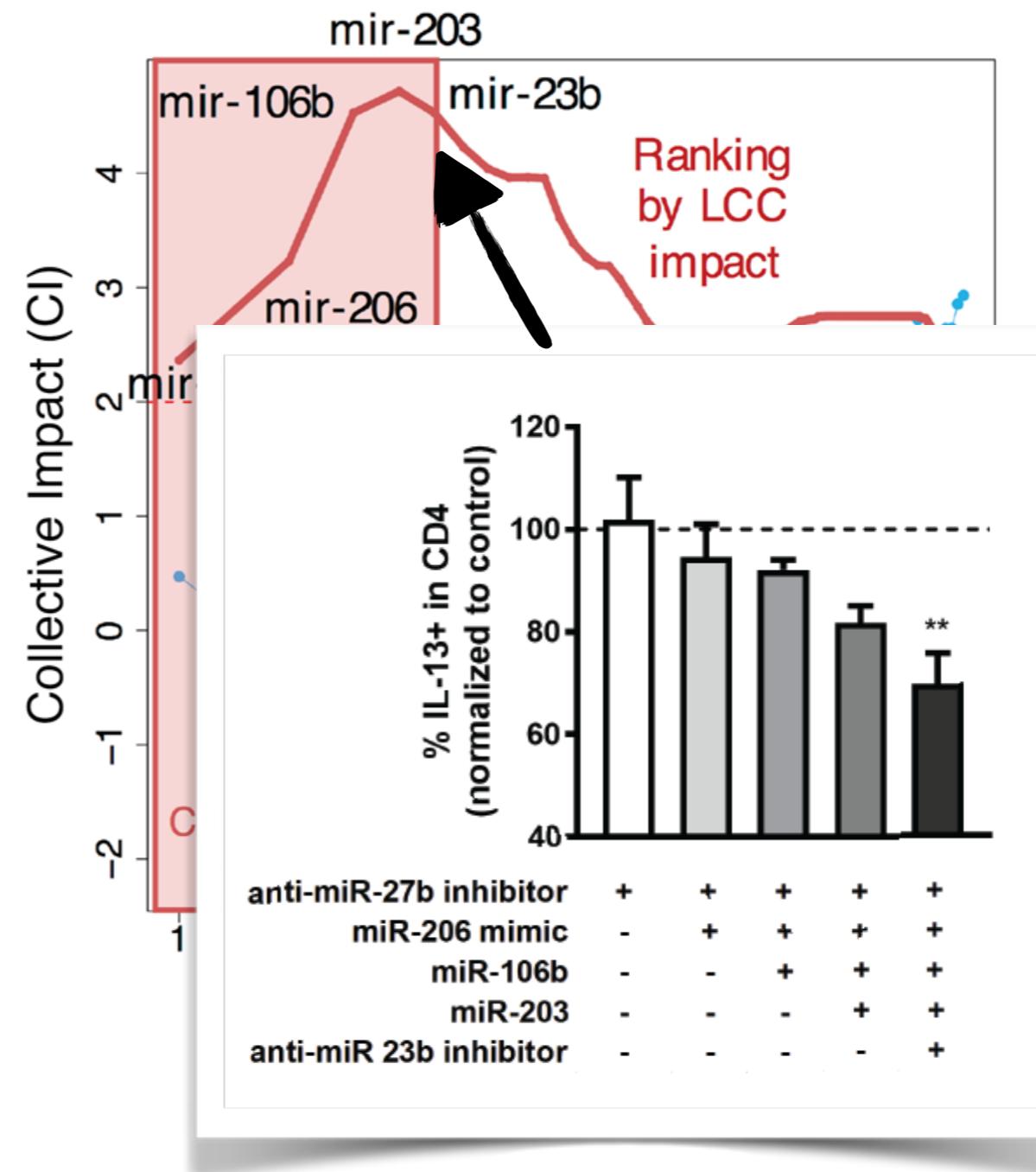
Jeong et al., "Lethality and centrality in protein networks" Nature 2001

EXAMPLE: MICRO-RNAs ATTACK

Network impact of miRNAs



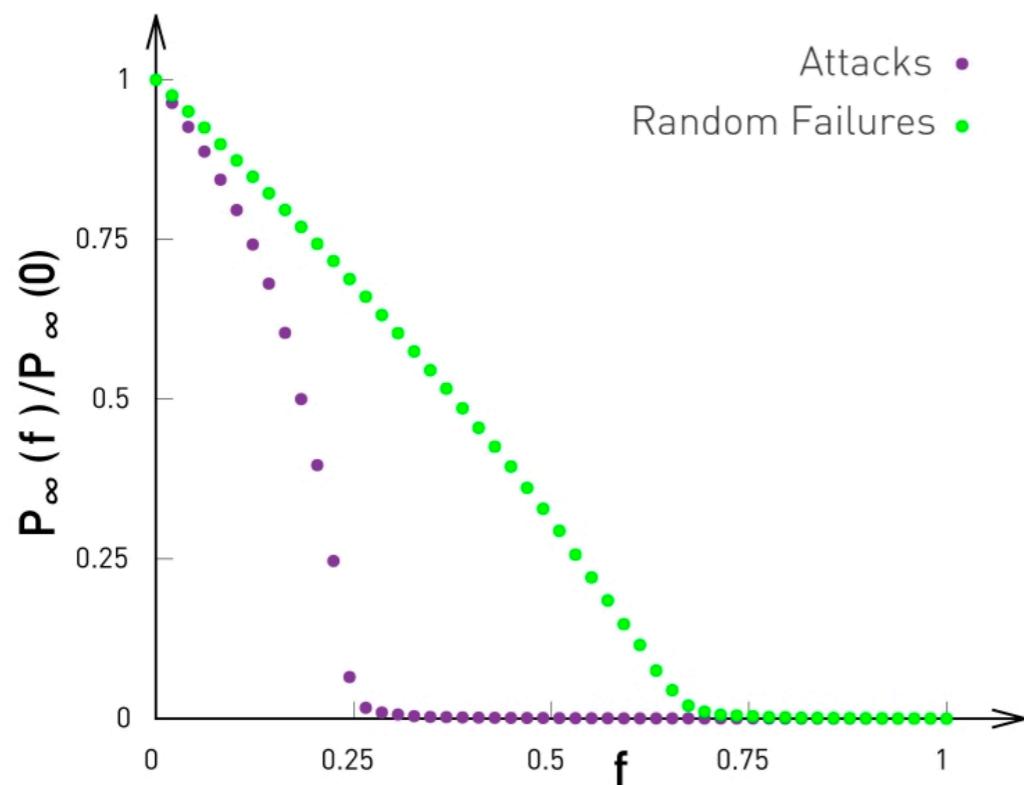
Kilic et al, JCI insight 2018



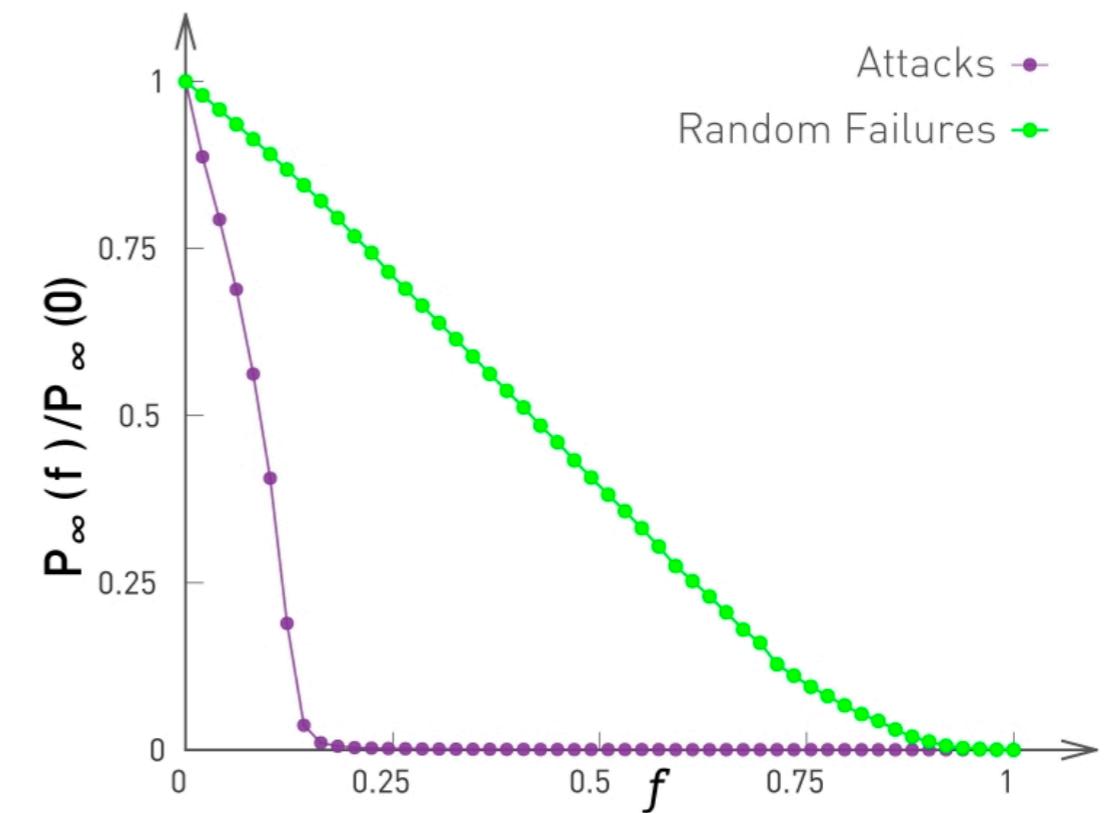
→ Network impact predicts functional miRNAs

Exercise: using python and networkx, reproduce the plots below.
 Display the network at 3 stages (initial, middle, end)
 file: resources Python /robustness_network.ipynb

Erdös-Renyi network



Barabasi-Albert network



Parameters

- $N = 10,000$
- $\langle k \rangle = 3$

Parameters

- degree exponent $\gamma = 2.5$
- $k_{\min} = 2$
- $N = 10,000$