

# How does the Web look like?

CS 322: (Social and Information) Network Analysis  
Jure Leskovec  
Stanford University



# Information networks

- Nodes are pieces of information
- Links join the pieces that are related to each other in some fashion
- **Examples:**
  - World Wide Web
  - Citation Networks
  - References in an encyclopedia
  - Internet
  - Wireless communication

# Before the Web was born

- Memex [*As we may think* by V. Bush, 1945]
  - Observed that information in books is highly linear
  - Our associative memory represents semantic network
  - **Memex** is a system that mimics such semantic networks with knowledge and links between the concepts

# Before the Web was born

## Gopher:

- Similar to Web:
  - links
- But:
  - No graphics
  - Imposes much stronger hierarchy – web of menu items – like a file system
  - Terminal based
  - “go for” (url)

```
Home Gopher server: gopherproject.org

THE GOPHER PROJECT
-----

Welcome to GOPHER! Gopher is a slim, powerful, and
fast way to present information in a hierarchial catalog online.
Gopher actually predates the Web -- although most web browsers
make excellent gopher browsers too.

Good places to start are the "Why Gopher?" and "Using Gopher"
areas!

--> [12] *** GOPHER TURNS 10 / GOPHER 3.0 (FurryTerror) RELEASED ***
[13] *** GOPHER TURNS 10 ..R 3.0 (FurryTerror) RELEASED *** [html] <HTML>
[14] A Brief Introduction to Gopherspace
[15] Clients, Servers, and Downloads/
[16] Home Gopher at UMN (a good place to browse)/
[17] Home Gopher at UMN [alternate]/
[18] Mailing List
[19] Mailing List Archives/
[20] Major Gopher Servers/
[21] Screenshots/

Press [h] for Help, [q] to Quit                                     Page: 1/2
```

# 1990: The Web is born

- Created by **Tim Berners-Lee** in 1990 in CERN:
  - a large hypertext database with typed links



# Christmas 1990: First Web page

## World Wide Web

The WorldWideWeb (W3) is a wide-area [hypermedia](#) information retrieval initiative aiming to give universal access to a large universe of documents.

Everything there is online about W3 is linked directly or indirectly to this document, including an [executive summary](#) of the project, [Mailing lists](#) , [Policy](#) , November's [W3 news](#) , [Frequently Asked Questions](#) .

### [What's out there?](#)

Pointers to the world's online information, [subjects](#) , [W3 servers](#), etc.

### [Help](#)

on the browser you are using

### [Software Products](#)

A list of W3 project components and their current state. (e.g. [Line Mode](#) ,X11 [Viola](#) , [NeXTStep](#) , [Servers](#) , [Tools](#) , [Mail robot](#) , [Library](#) )

### [Technical](#)

Details of protocols, formats, program internals etc

### [Bibliography](#)

Paper documentation on W3 and references.

### [People](#)

A list of some people involved in the project.

### [History](#)

A summary of the history of the project.

### [How can I help ?](#)

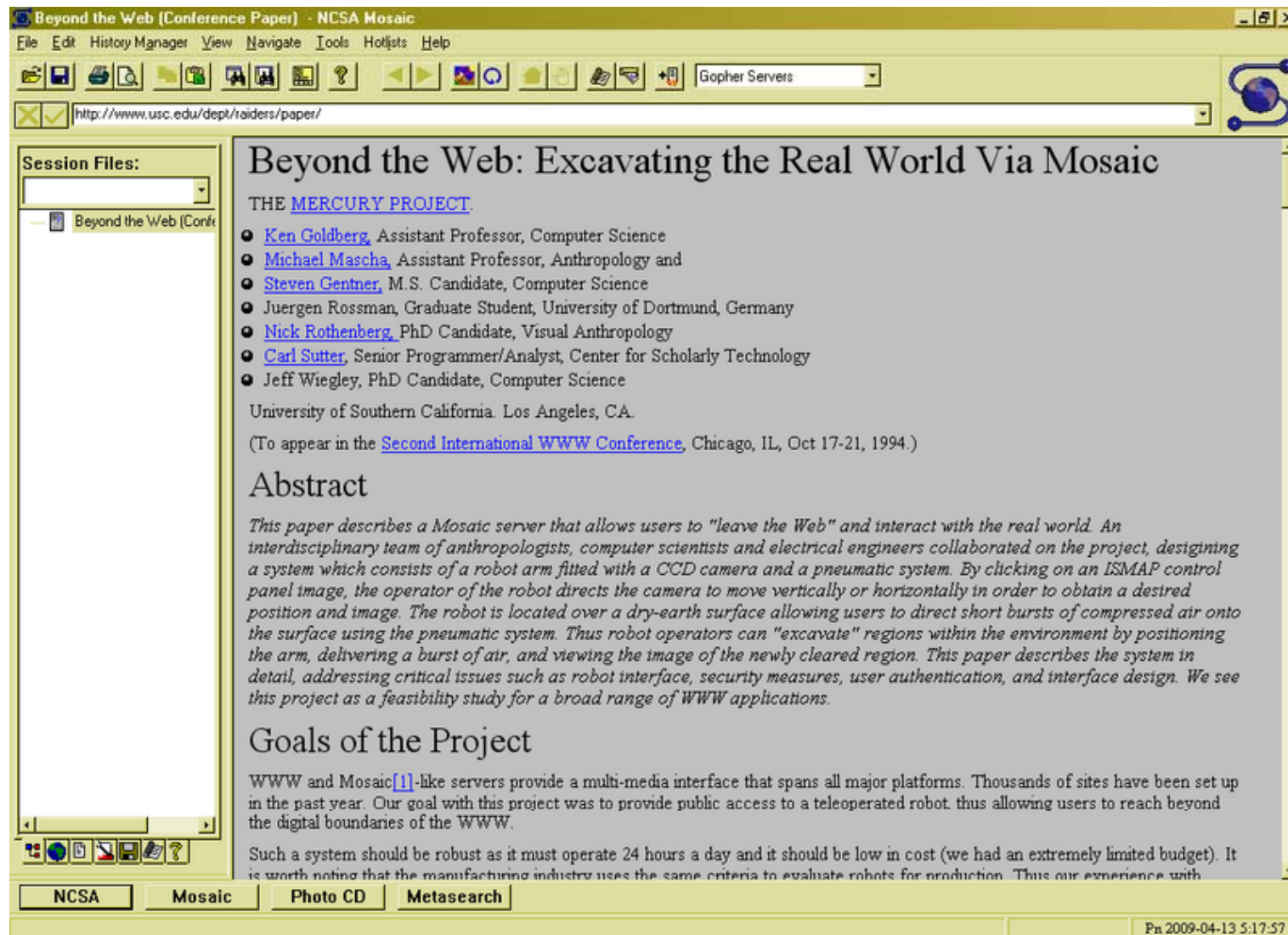
If you would like to support the web..

### [Getting code](#)

Getting the code by [anonymous FTP](#) , etc.

<http://www.w3.org/History/19921103-hypertext/hypertext/WWW/TheProject.html>

# 1993: First real browser – Mosaic



# Brief history of the Web

- 1989:
  - HTTP, HTML
- 1990, Christmas:
  - First web server at CERN
- 1992-95: Growth
  - 1993: first Unix graphical web browser – Mosaic
  - 1994: first WWW conference, W3C is formed
- 1996-98:
  - Commercialization of the WWW
- 1999-2001:
  - Dot-com boom
- 2002:
  - Web is ubiquitous
    - Web 2.0,
    - User generated content (blogs, rss)
    - Semantic Web



- [Arts](#) -- [Humanities](#), [Photography](#), [Architecture](#)
- [Business and Economy \[Xtra!\]](#) -- [Directory](#)
- [Computers and Internet \[Xtra!\]](#) -- [Internet](#), [WWW](#), [Software](#), [Multimedia](#) ...

Yahoo, 1996



Amazon, 1995



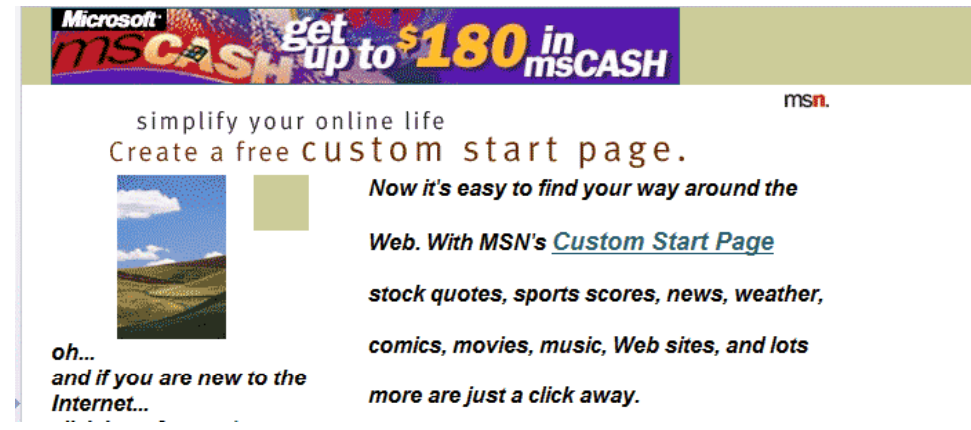
# Brief history of the Web



Microsoft, 1995



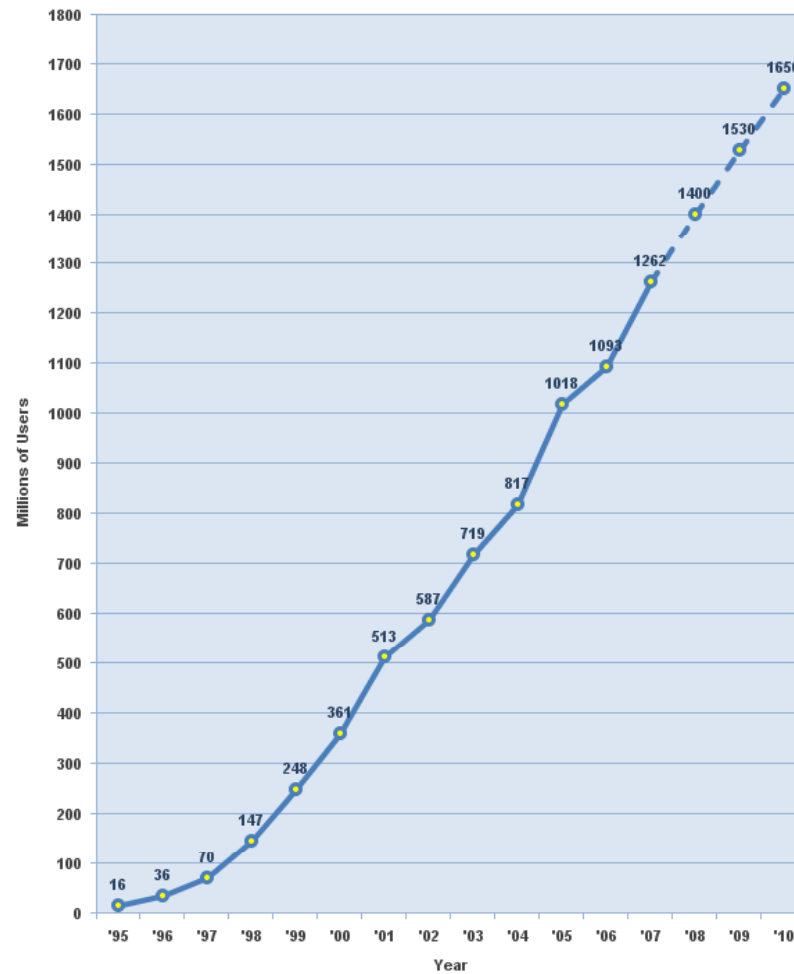
Google, 1998



MSN, 1996

# Growth of the Web: Users

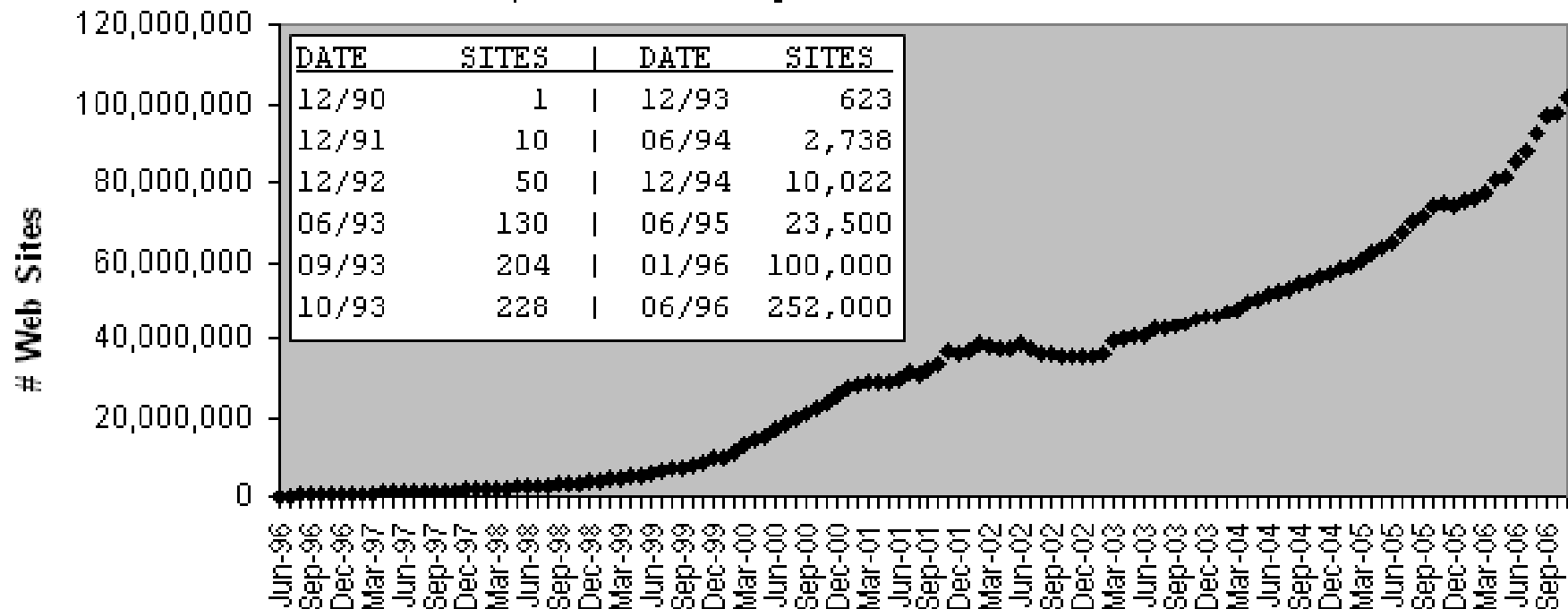
Internet Users in the World  
Growth 1995 - 2010



Source: [www.internetworldstats.com](http://www.internetworldstats.com) - January, 2008  
Copyright © 2008, Miniwatts Marketing Group

# Growth of the Web: Websites

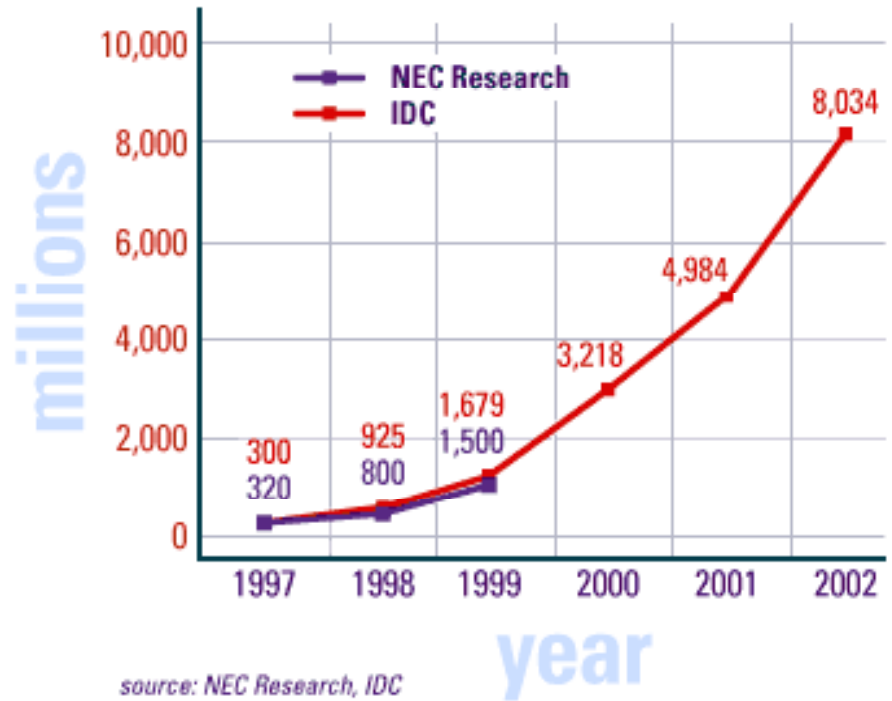
Hobbes' Internet Timeline Copyright ©2006 Robert H Zakon  
<http://www.zakon.org/robert/internet/timeline/>



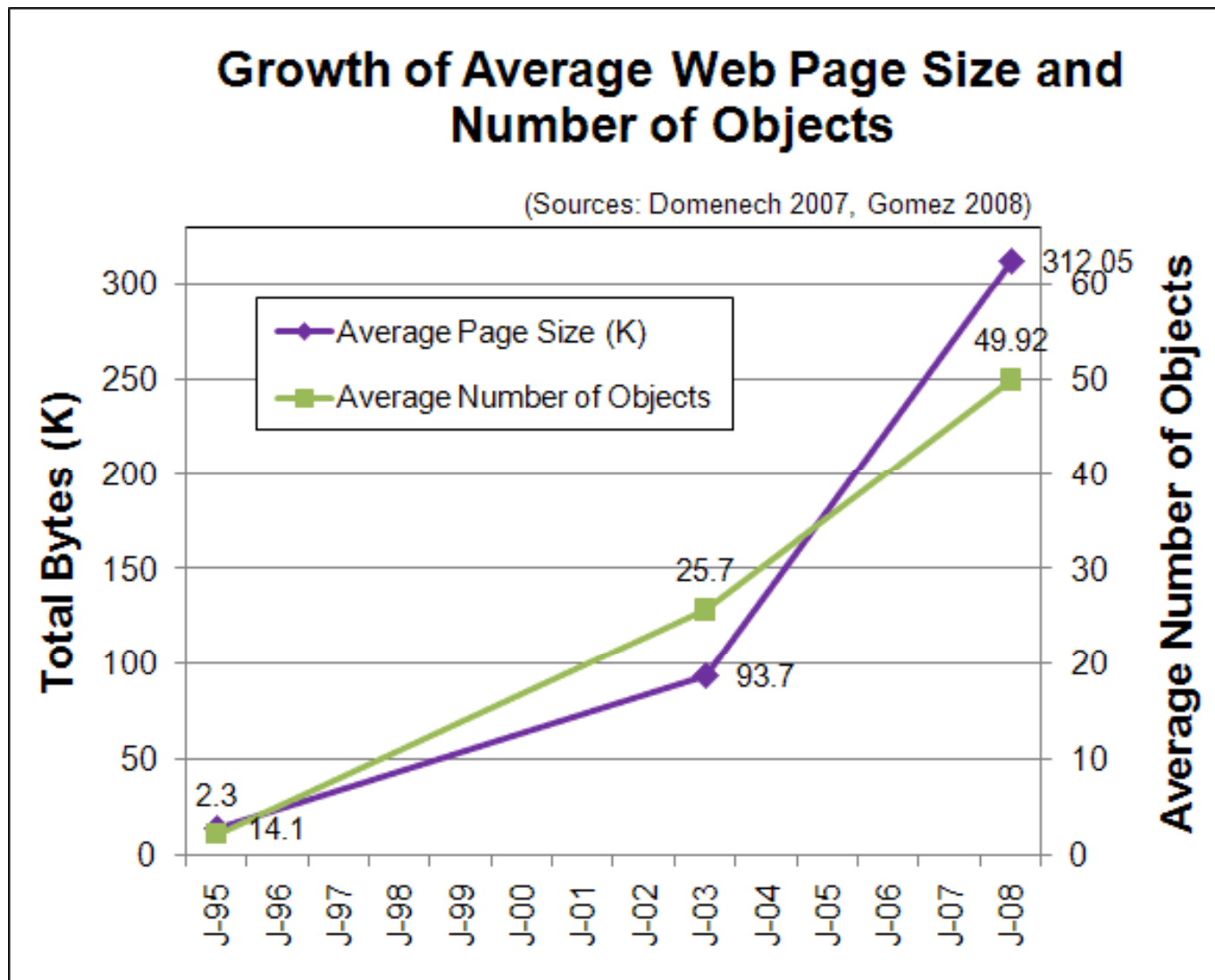
- Estimated number of websites (web servers)
- August 2009: 225 million hostnames

# Growth of the Web: Pages

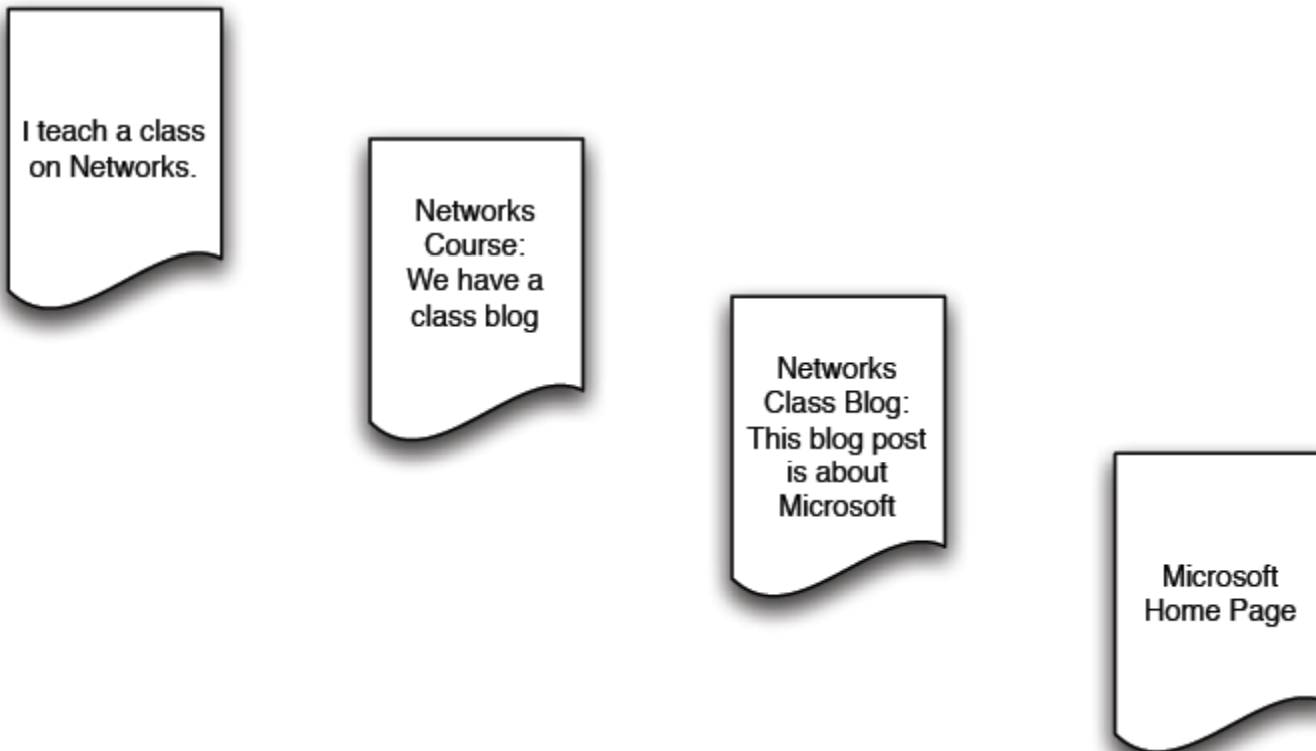
- Yahoo, 2005:
  - indexed 19.2 billion pages
- Feb 2007:
  - Estimate 29 billion pages
- Not clear:
  - How do we discover webpages?
  - Hidden web (deep web)
  - What does a page mean:
    - dynamically generated pages
    - Time stamped URLs



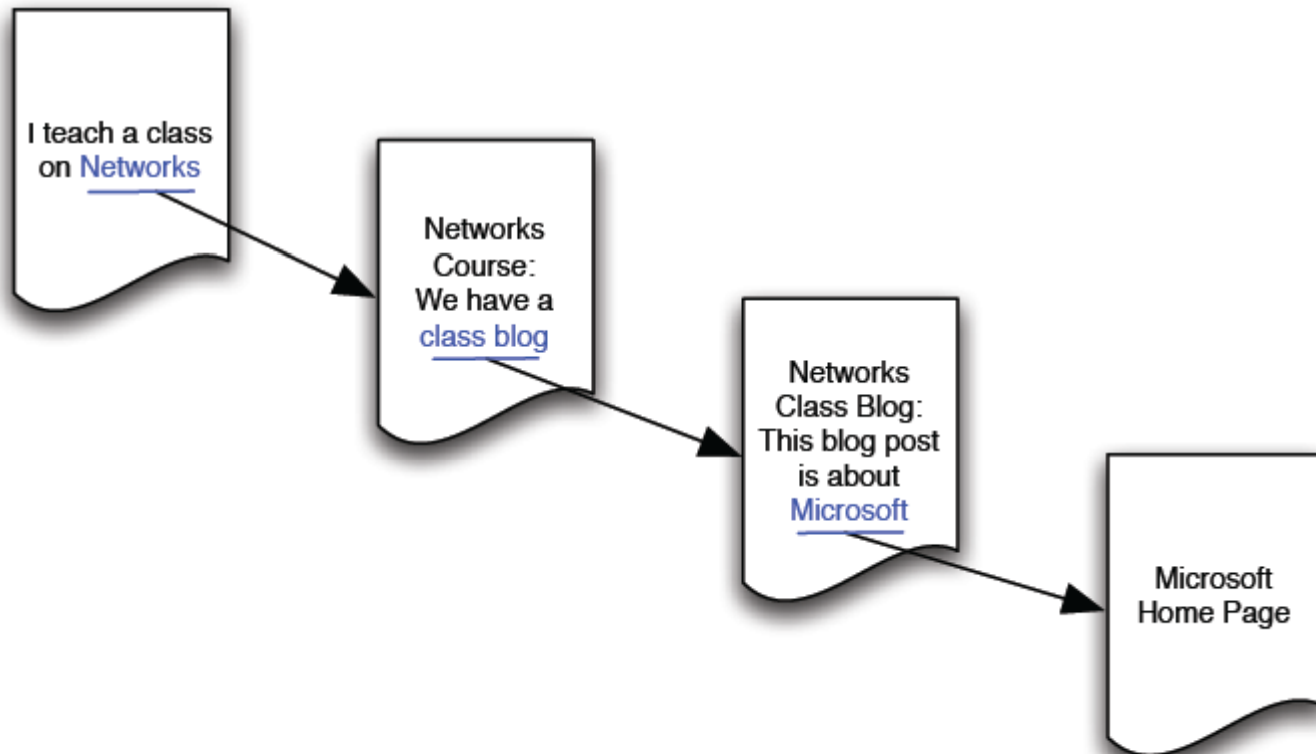
# Growth of the Web: Page size



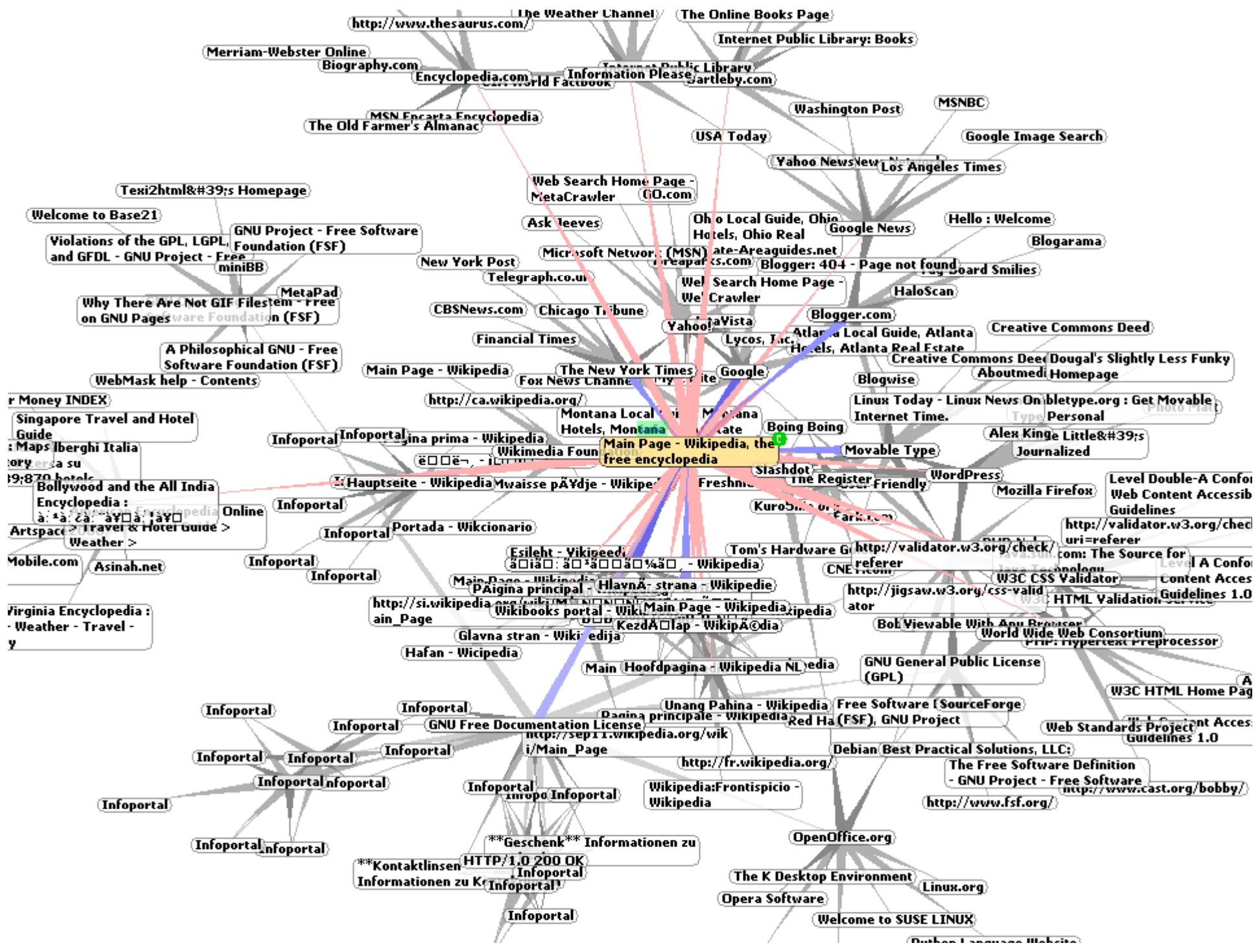
# Web as a Graph



# Web as a Graph

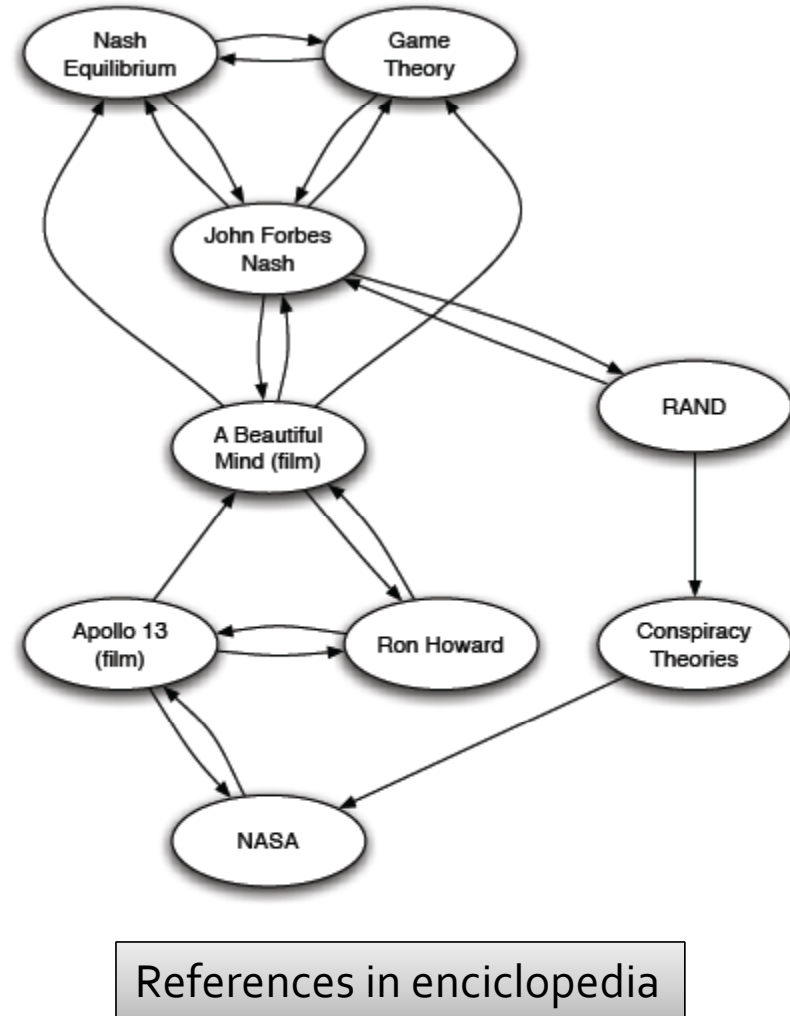
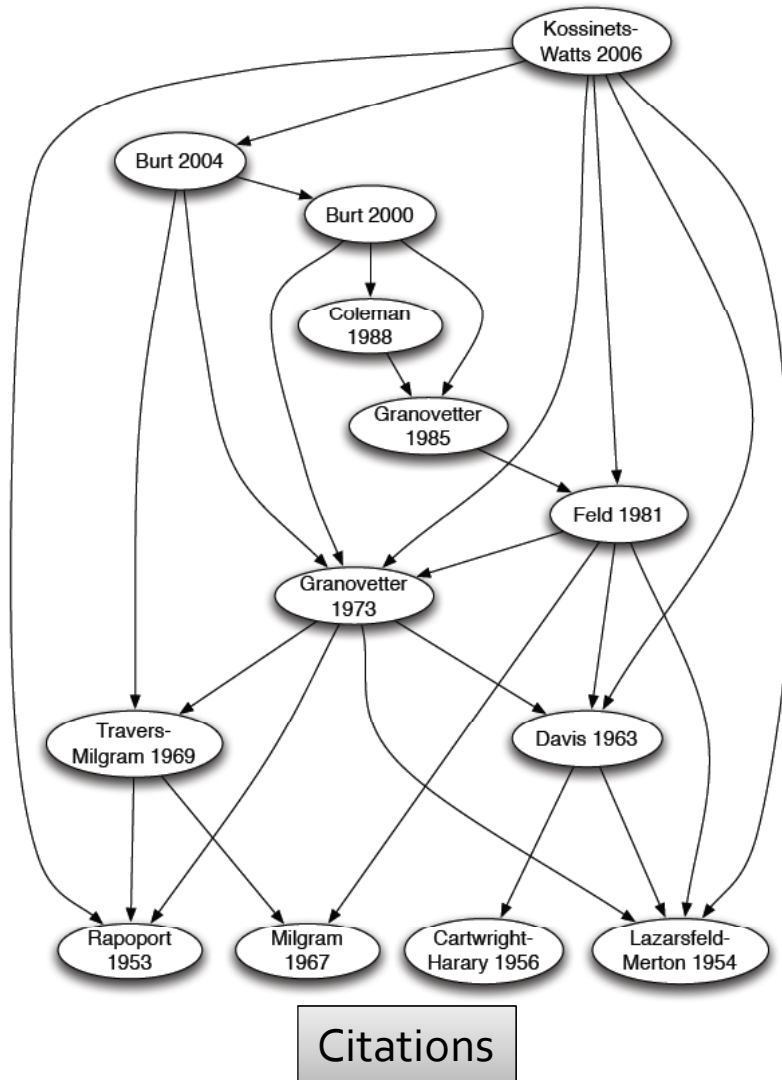


- In early days of the Web links were **navigational**
- Today many links are **transactional**





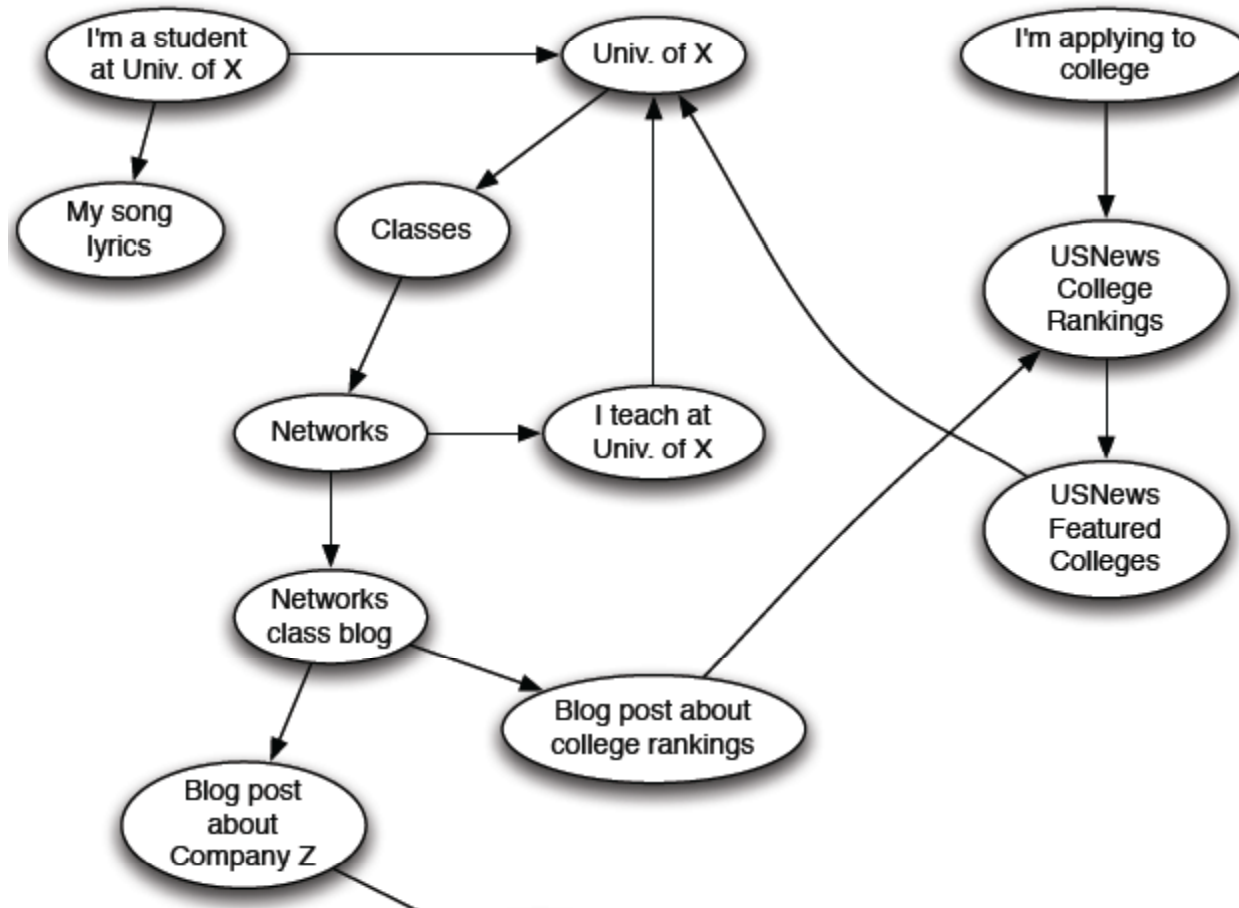
# Other information networks



# How does the Web look like?

- How is information on the Web organized?
- How is the Web linked?
- What is the “map” of the Web?

# Web as a directed graph



# Directed graphs

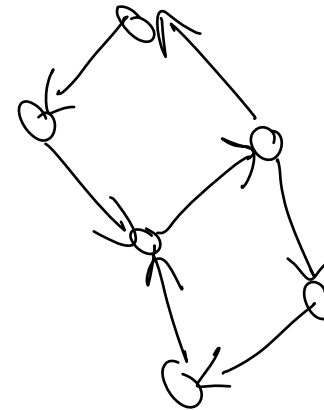
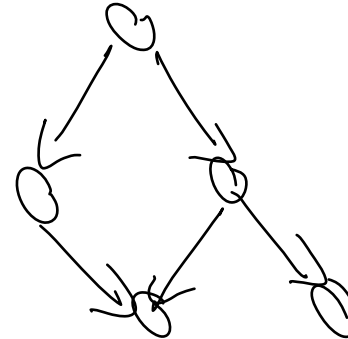
- Two types of directed graphs:

- **DAG – directed acyclic graph:**

- Has no cycles: if  $u$  can reach  $v$ , then  $v$  can not reach  $u$

- **Strongly connected:**

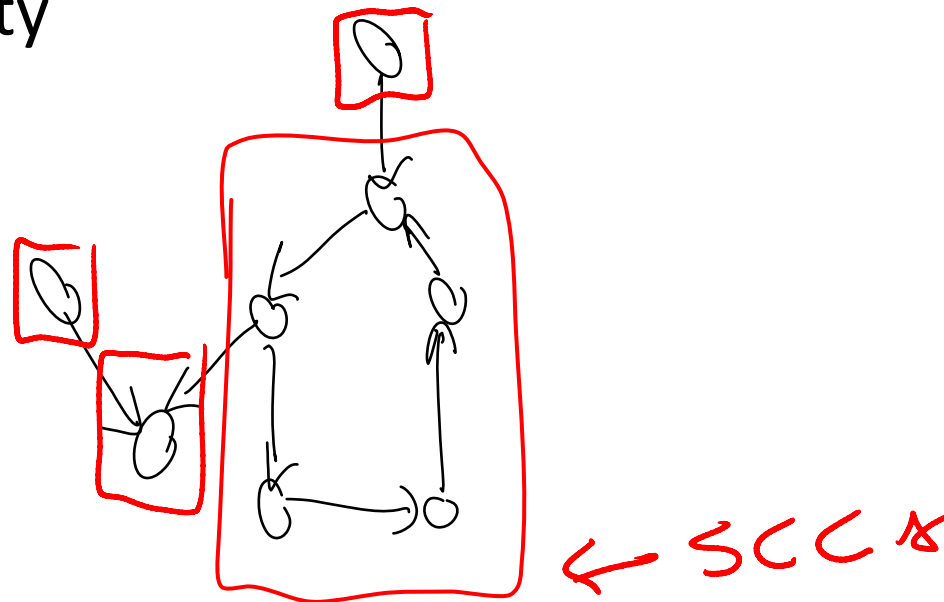
- Any node can reach any node via a directed path



- Any directed graph can be expressed in terms of these two types

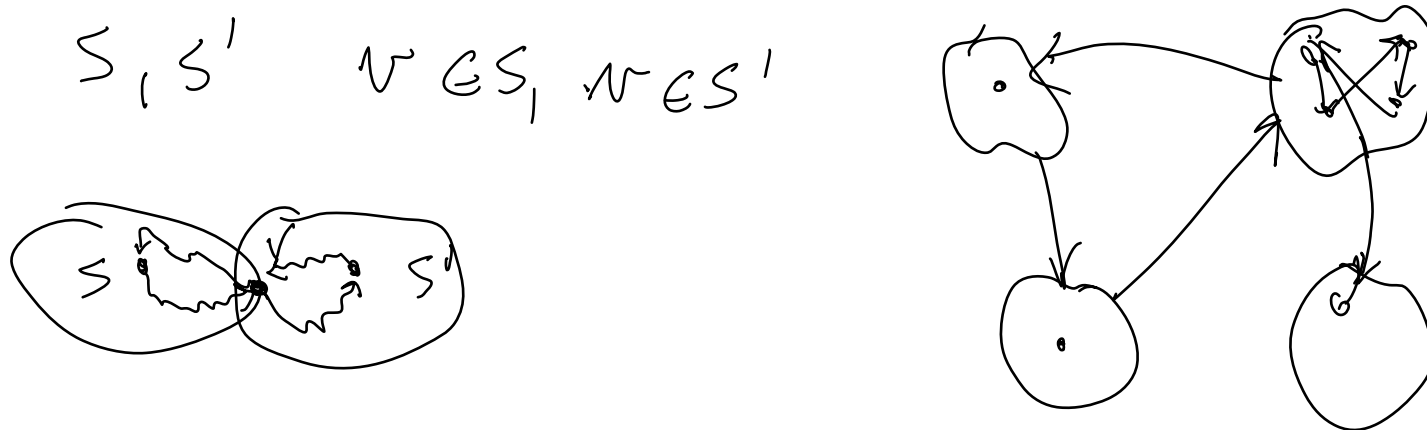
# Strongly connected component

- Strongly connected component (SCC) is a set of nodes  $S$  so that:
  - Every pair of nodes in  $S$  can reach each other
  - There is no larger set containing  $S$  with this property



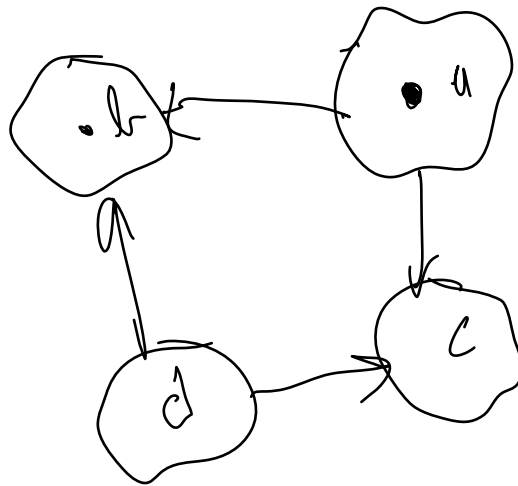
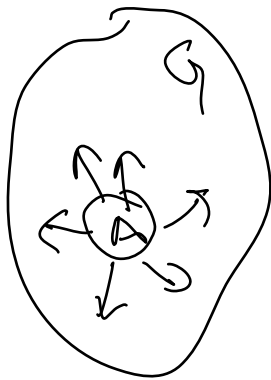
# Strongly connected component

- Fact: Every directed graph is a DAG on its SCCs.
  - SCCs partition the nodes of  $G$  (each node in exactly one SCC)
  - If we build a graph  $G'$  whose nodes are SCCs, and edge between nodes  $G'$  if there is an edge between corresponding SCCs in  $G$ , then  $G'$  is a DAG



# Example

- Given a directed graph  $G$
- Find the smallest set of node  $A$  so that every node in  $G$  is reachable from at least one node in  $A$ .



$u \in A$

# Graph structure of the Web

- Take a large snapshot of the web and try to understand how its SCCs “fit” as a DAG.

- Computational issue:**

- Say want to find SCC containing specific node  $v$ ?

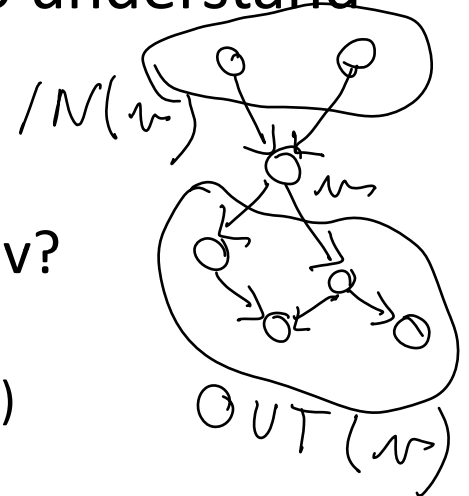
- Observation:**

- Out( $v$ ) ... nodes that can be reached from  $v$  (BFS out)
- SCC containing  $v$  is:

$$= \text{OUT}(v) \cap \text{IN}(v)$$

$$= \text{OUT}(v, G) \cap \text{OUT}(v, \bar{G})$$

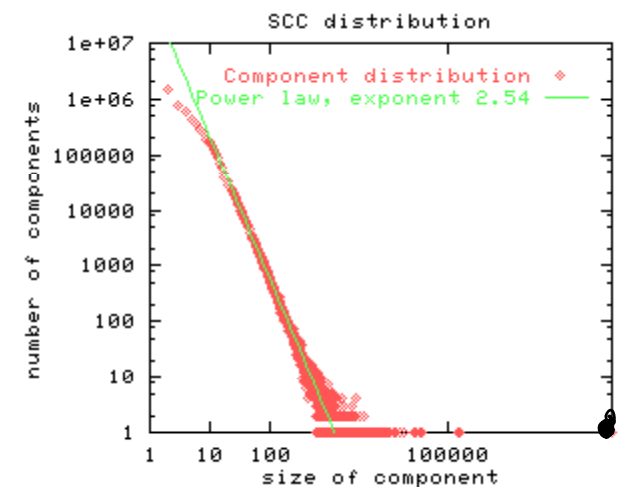
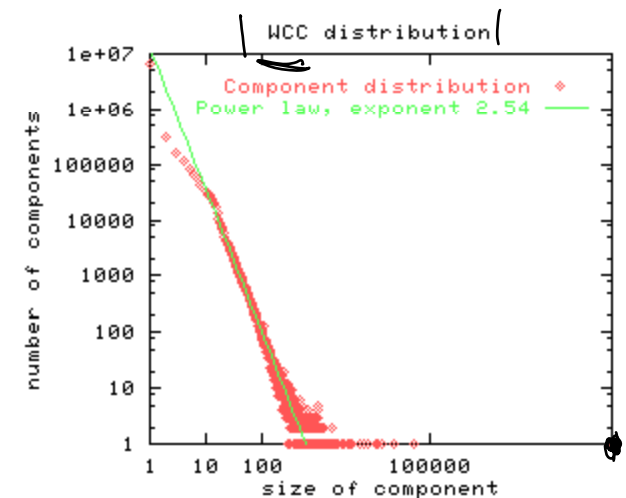
$\bar{G}$  ...  $G$  with edge directions flipped



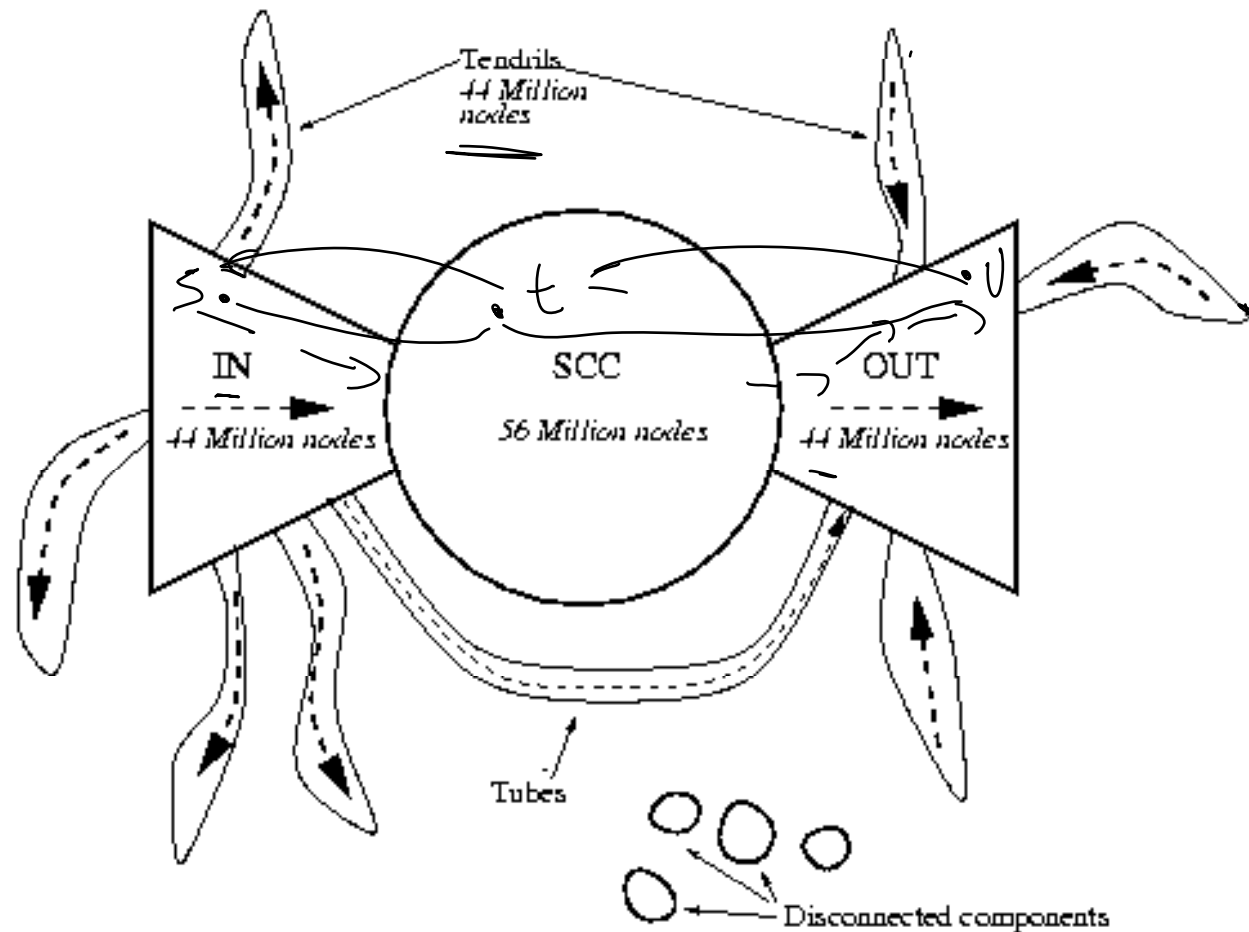


# Graph structure of the Web

- There is a giant SCC
- There will not be 2 giant SCCs:
  - Just takes 1 page from each to link to one in other – if the components have millions of pages the likelihood of this is large
- Broder et al., 2000:
  - Weakly connected component: 90% of the nodes

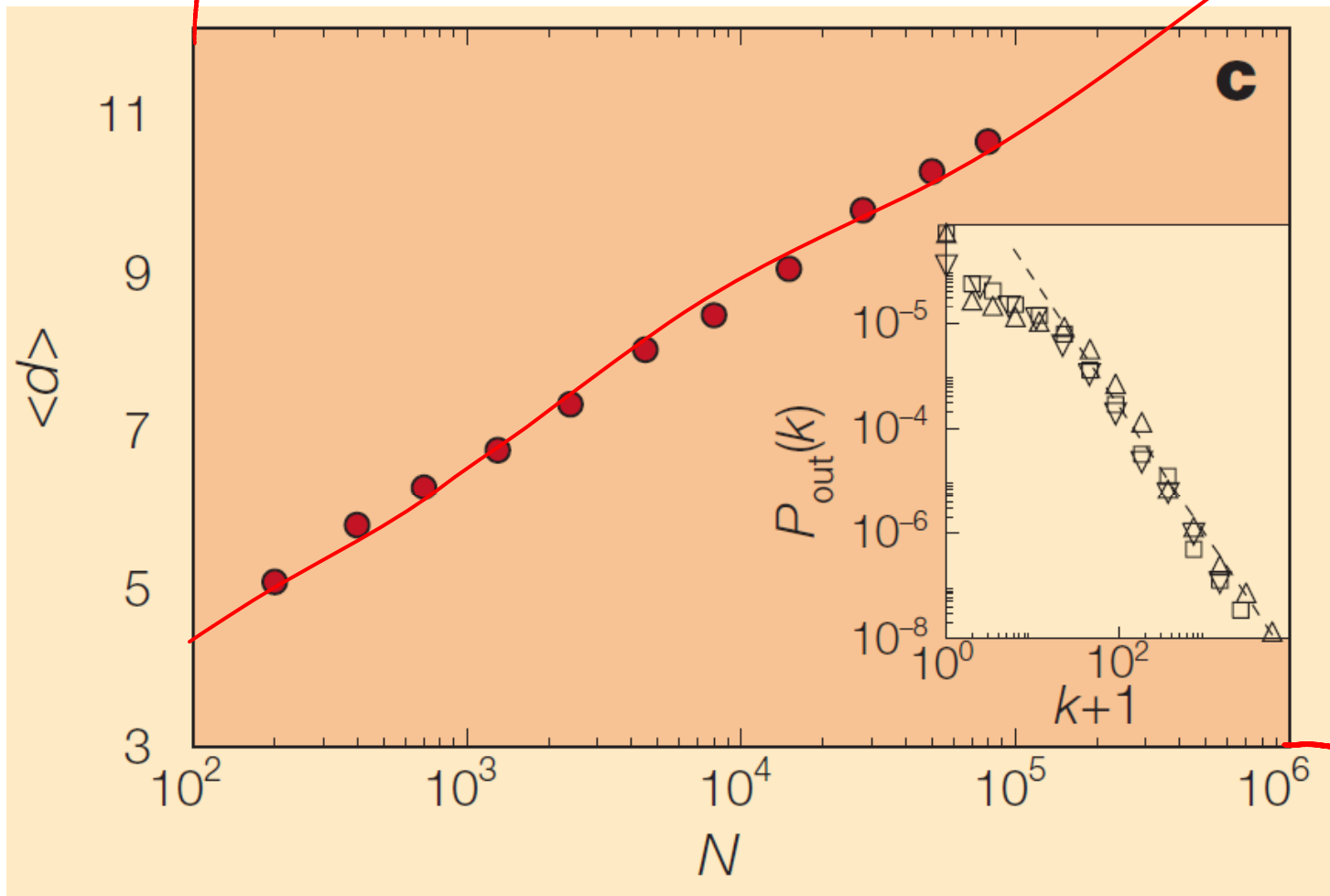


# Bow-tie structure of the Web



- 250 million webpages, 1.5 billion links [Altavista]

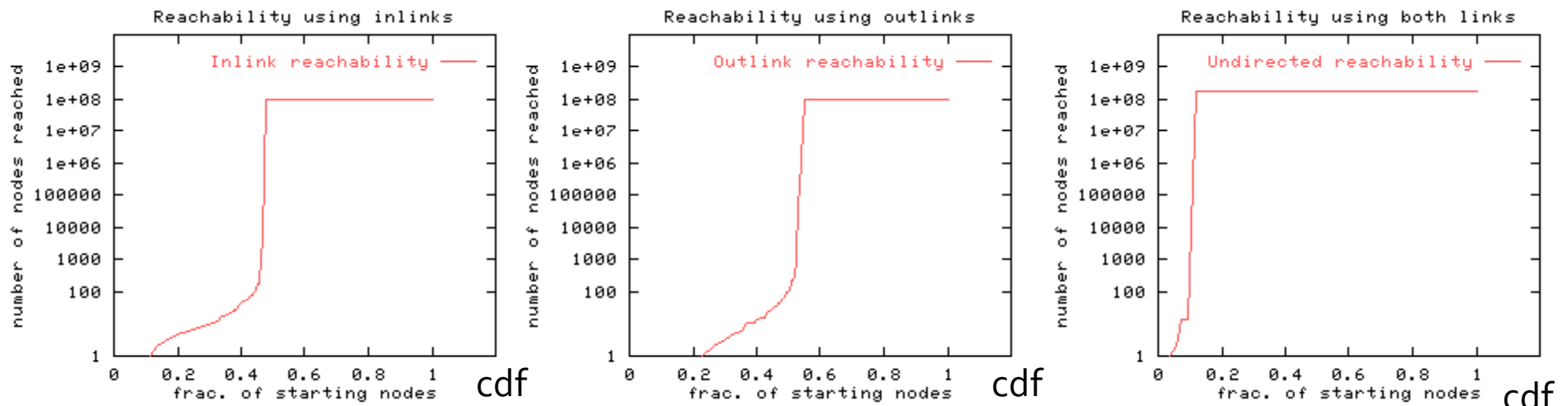
# Diameter of the Web



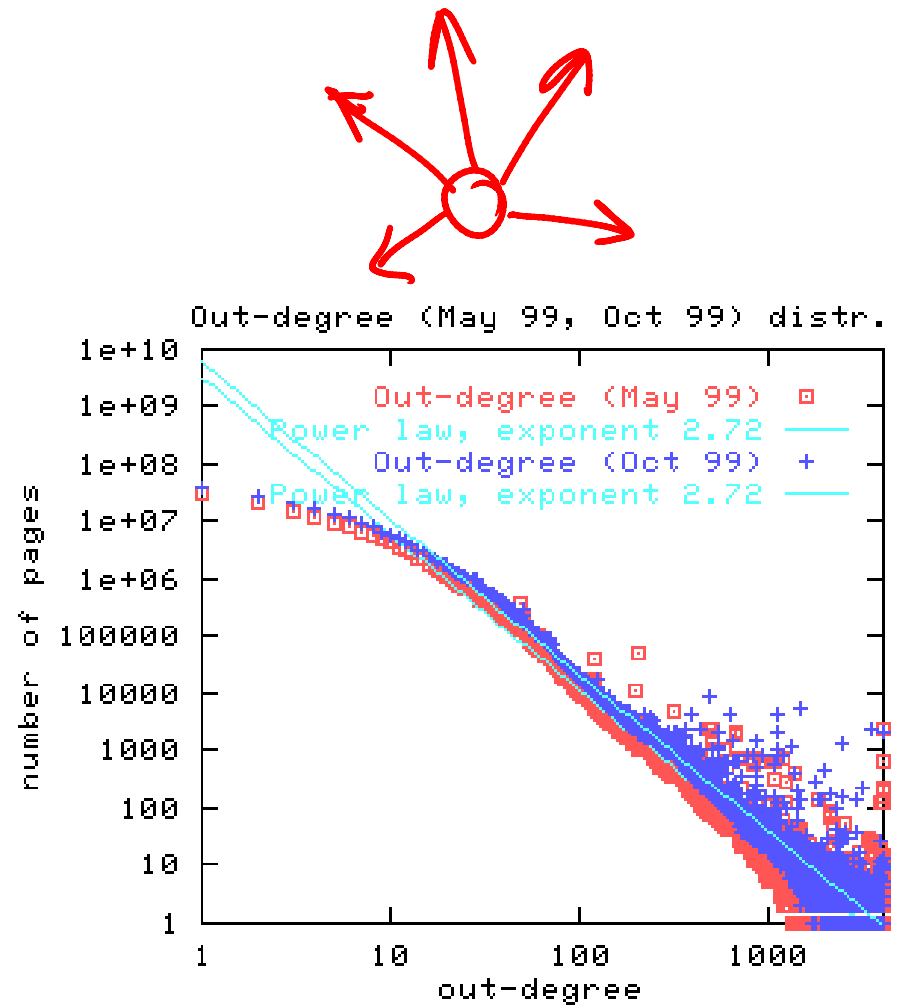
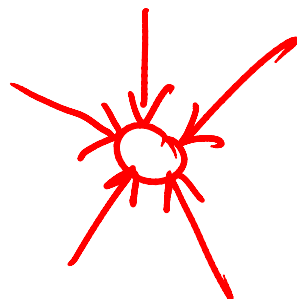
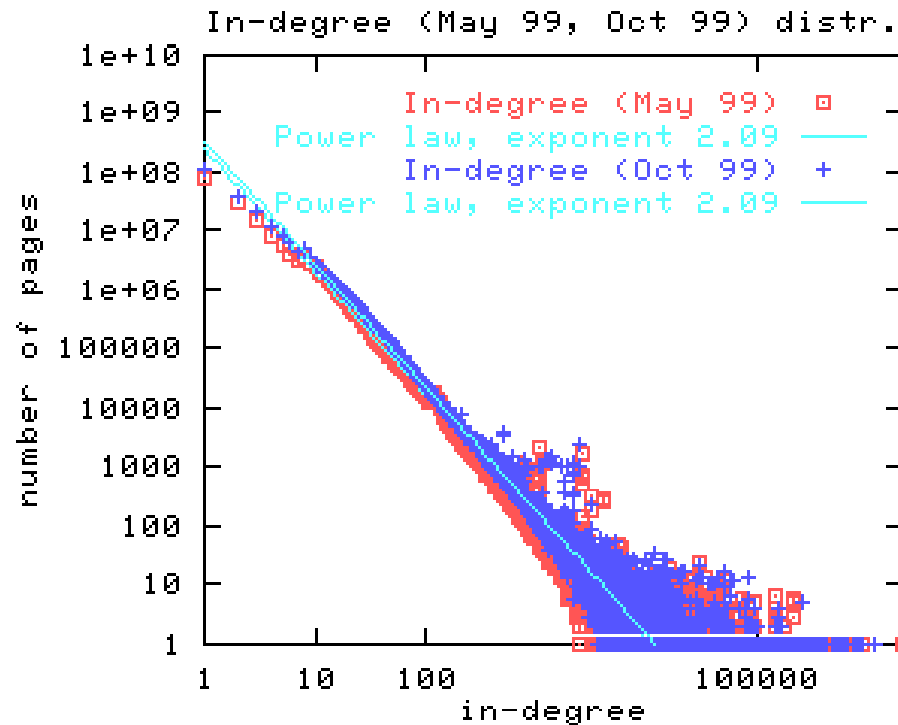
- Diameter (average directed shortest path length) is 19 (in 1999)

# Diameter of the Web

- Average distance: 75% of time there is no directed path
  - Follow in-links (directed): 16.12
  - Follow out-links (directed): 16.18
  - Undirected: 6.83
- Diameter of SCC (directed): At least 28

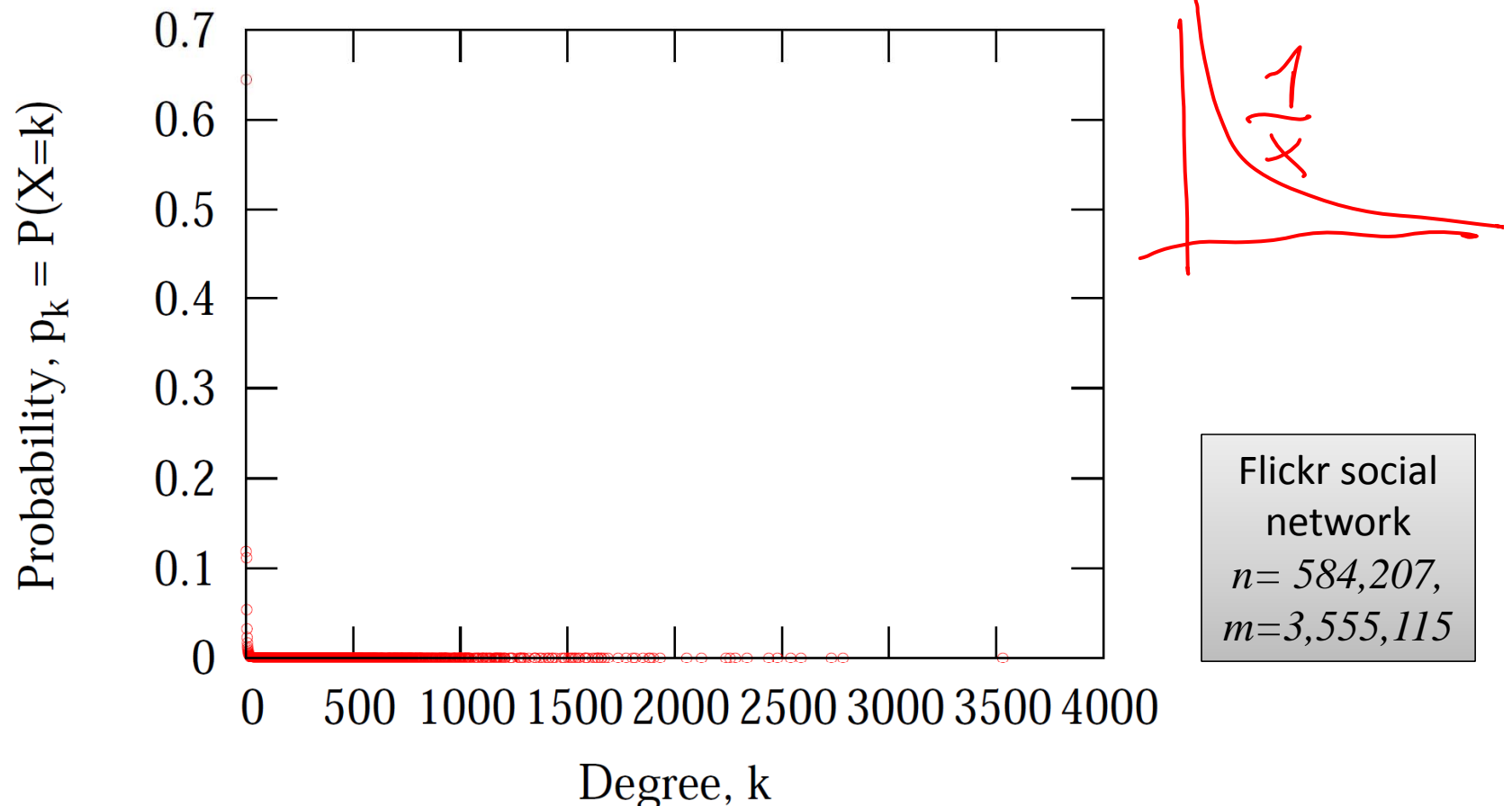


# Degree distribution on the Web



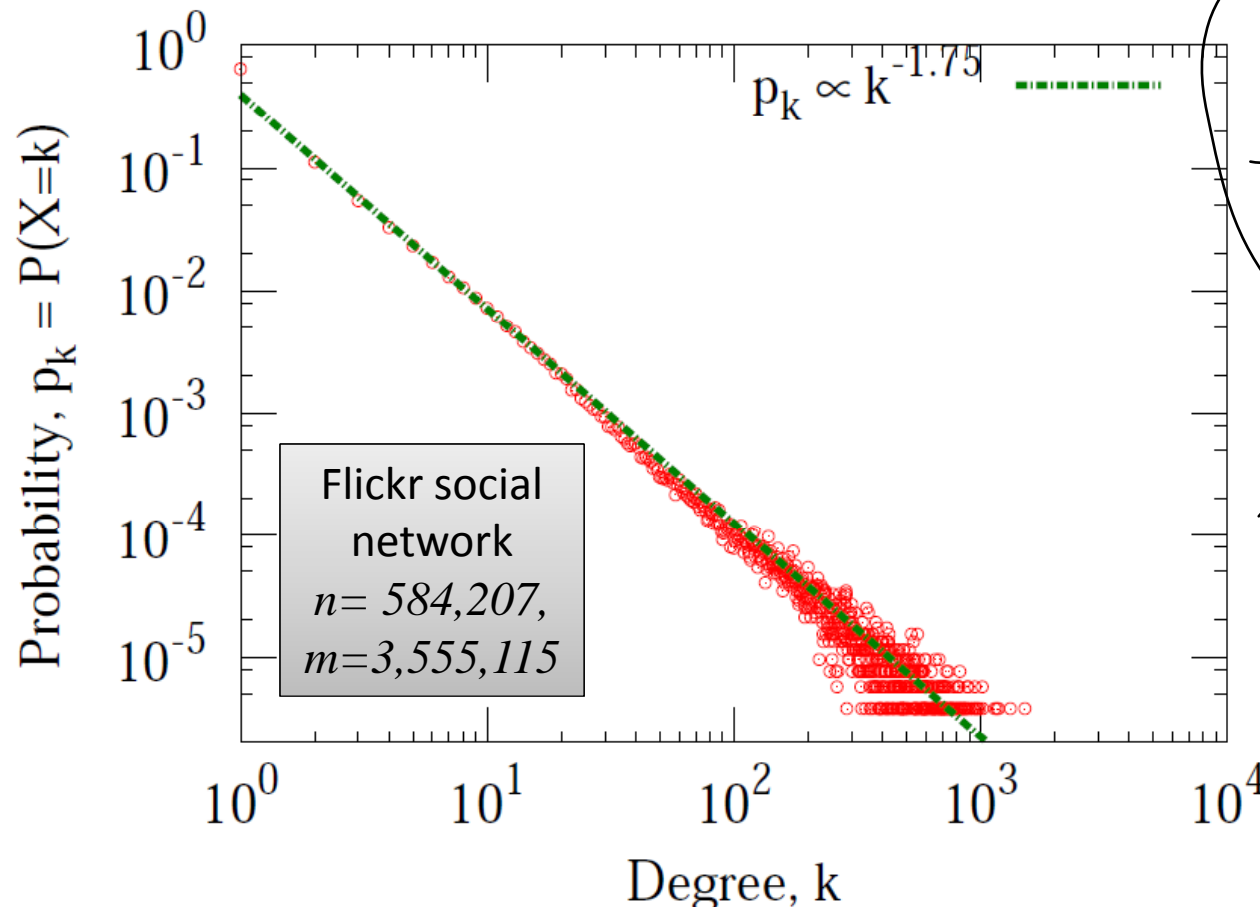
# Degrees in real networks

- Take real network plot a histogram of  $p_k$  vs.  $k$



# Degrees in real networks (2)

- Plot the same data on *log-log* axis:  $f(k) \sim c \cdot k^{-\alpha}$



$$f(k) \sim c \cdot k^{-\alpha}$$


---


$$\log f(k) = \log c$$

$$+ -\alpha \log k$$


---

slope  $-\alpha$

---


$$\log f(k)$$

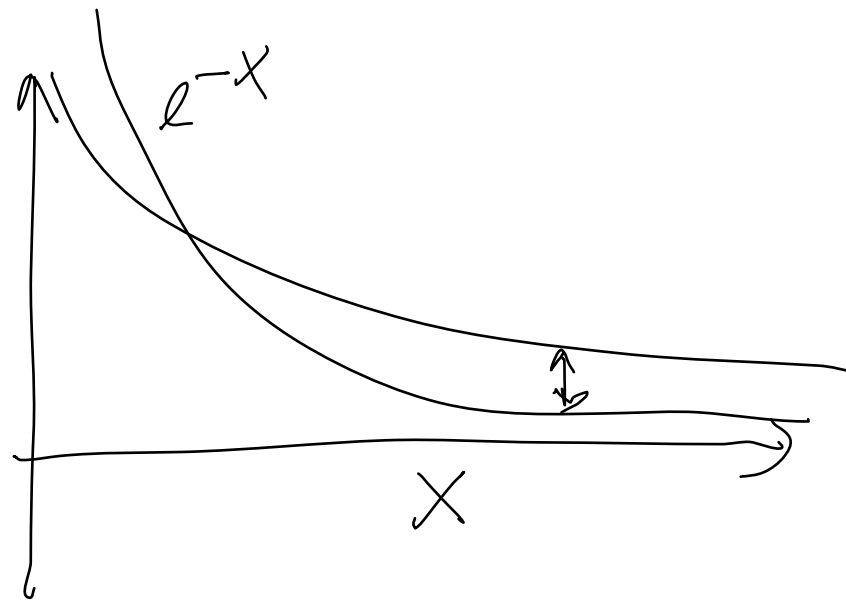

---


$$\log k$$

# Heavy tails

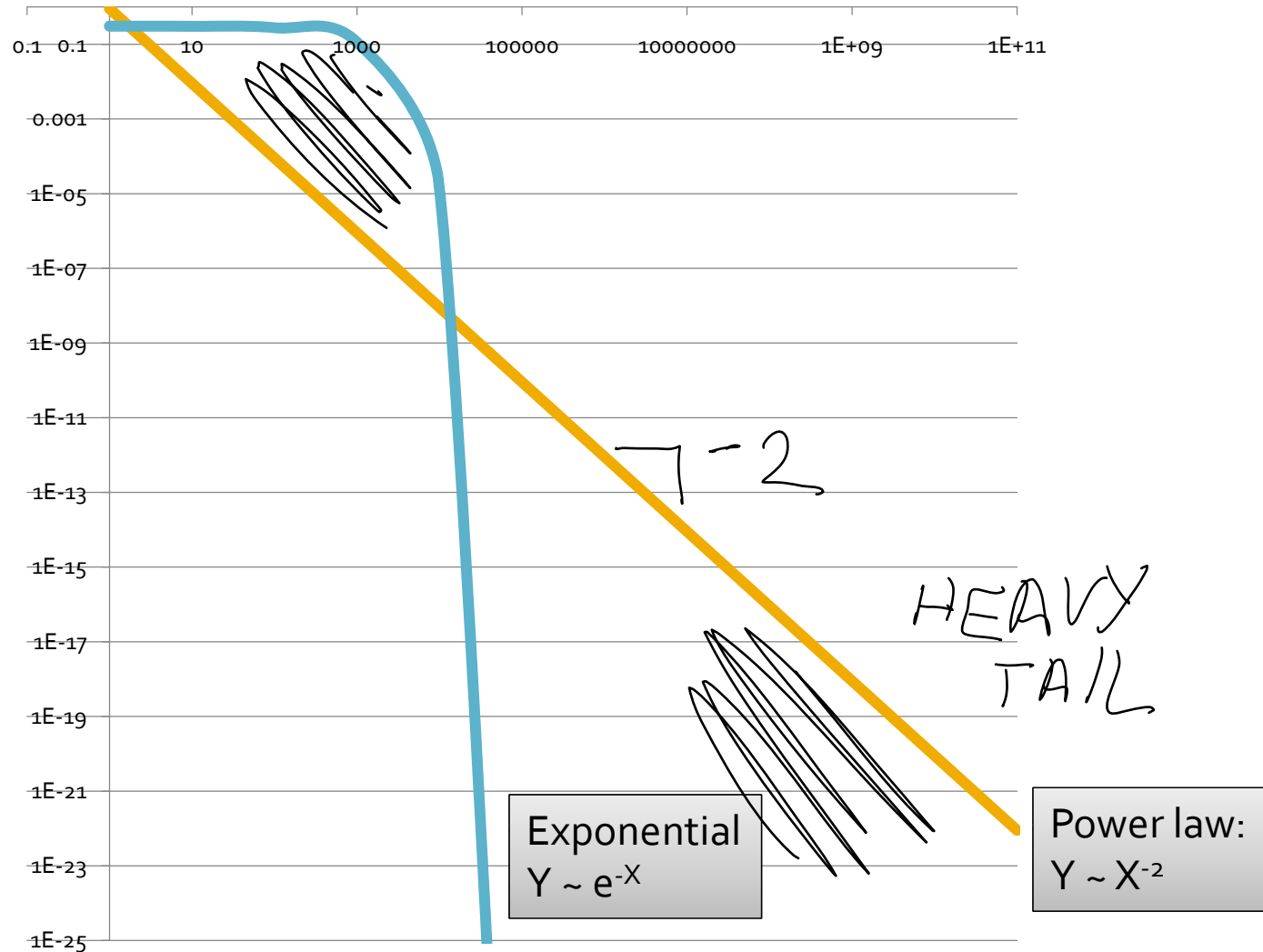
- Degrees are heavily skewed:  
Distribution is **heavy tailed**:

$$\lim_{x \rightarrow \infty} \frac{\Pr(X > x)}{e^{-\epsilon x}} = \infty$$





# Exponential tail vs. Power-law tail



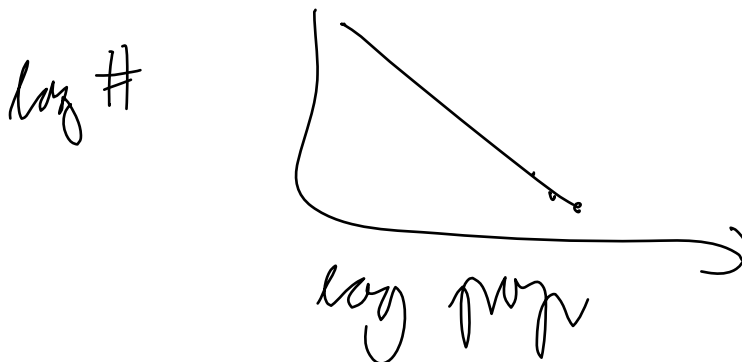
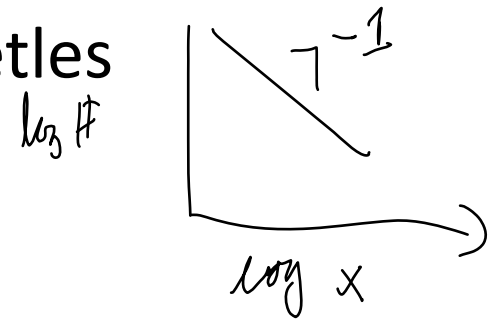
# Heavy tails

- Various names, kinds and forms:
  - Long tail, Heavy tail, Zipf's law, Pareto's law, 80-20 rule

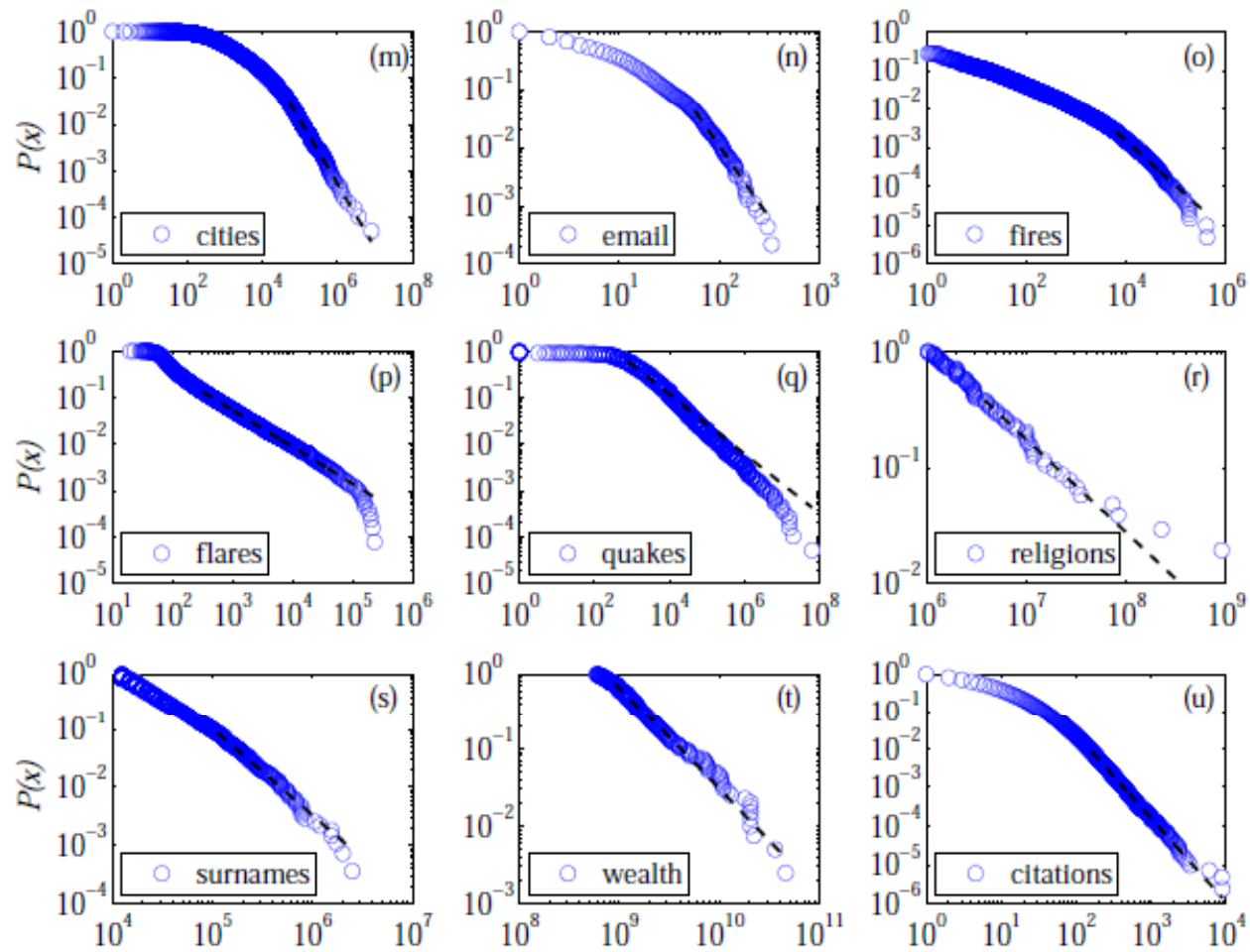
power law	$p(x) \propto \frac{x^{-\alpha}}{\quad}$
power law with cutoff	$x^{-\alpha} e^{-\lambda x}$
stretched exponential	$x^{\beta-1} e^{-\lambda x^{\beta}}$
log-normal	$\frac{1}{x} \exp \left[ -\frac{(\ln x - \mu)^2}{2\sigma^2} \right]$

# Power-laws are everywhere

- In social systems – lots of power laws:
  - Pareto, 1897 – Wealth distribution
  - Lotka 1926 – Scientific output
  - Yule 1920s -- lots of species of beetles
  - Zipf 1940s – word frequency
  - Simon 1950s – city populations



# Power-laws are everywhere



Many other quantities follow heavy-tailed distributions

# The long tail

## ANATOMY OF THE LONG TAIL

Online services carry far more inventory than traditional retailers. Rhapsody, for example, offers 19 times as many songs as Wal-Mart's stock of 39,000 tunes. The appetite for Rhapsody's more obscure tunes (charted below in yellow) makes up the so-called Long Tail. Meanwhile, even as consumers flock to mainstream books, music, and films (right), there is real demand for niche fare found only online.



Sources: Erik Brynjolfsson and Jeffrey Hu, MIT, and Michael Smith, Carnegie Mellon; Barnes & Noble; Netflix; RealNetworks

# Not everyone likes power-laws 😊



CMU students protesting at the G20 meeting in Pittsburgh in Sept 2009