



# Bibliometrics of the Web

**Ray R. Larson**  
Professor

**School of Information Management and Systems**  
**University of California, Berkeley**

# Overview



- Web Characteristics
- Bibliometrics/Webometrics
- Uses of Bibliometrics (link analysis)
- New Approaches and New Uses

# Overview



- **Web Characteristics**
- Bibliometrics/Webometrics
- Uses of Bibliometrics (link analysis)
- New Approaches and New Uses

# Web Characteristics



- The WWW differs radically from all previous information systems, lacking any central imposition of control and the ability for anyone or any group to develop sites (or a "web presence")
- The WWW is an "'ecological' information system" with "
  - self-organized clusters and communities
  - islands of order + oceans of topical diversity
  - local actions → global consequences (e.g. small-world phenomena)
  - '3D' = distributed + diversive + dynamic"

» Bjorneborn *Small-world link structures on the Web* 2002

# Web Size and Links



- Estimated 3.5-5 Billion "visible" web pages
  - Not including over 8Tb in the 'Deep Web'
- Estimated 35-50 Billion links
  - Dynamic link structures
    - Links are created, modified and destroyed by millions of individuals, institutions, companies, programs, etc.
    - No real typology for web links -- there is an infinity of possible reasons for linking behaviour including similarity and other topical relations, personal preferences or interests, navigational means
  - Dynamic adaptation and removal of links
    - may reflect changes in web constructors' knowledge, ideas, interests, contacts, etc.
    - may reflect cultural and social currents, e.g., diffusion of scientific ideas

# Overview



- Web Characteristics
- **Bibliometrics/Webometrics**
- Uses of Bibliometrics (link analysis)
- New Approaches and New Uses



- Bibliometrics is concerned with analysis of the characteristics of published literature including:
  - Studies of dispersion of literature on various topics
  - Statistical analyses of content types, references, etc.
  - Citation and co-citation studies within and across particular disciplines

# Bibliometrics and Webometrics



- Larson (1996) applied co-citation analysis to the WWW
- Webometrics (Almind and Ingwersen 1997) is the application of bibliometric methods to the WWW
- Webometrics is similar to traditional informetric and scientometric studies, applying common bibliometric methods.
  - counts and content analysis of web pages are like traditional publication analysis
  - counts and analyses of outgoing links from web pages (*outlinks*) and links pointing to web pages (*inlinks*) are very like reference and citation analyses
    - However circular links (where a page references itself or a part of itself) may not fulfil the same role in pre-web published literature



# Citation Analysis



- **Citation analysis** was developed in information science as a tool to identify core sets of articles, authors, or journals of particular fields of study.
- **Cocitation analysis** has been used to map the topical relatedness of clusters of authors, journals or articles (see H. White 1981)
- It can provide a mapping of the intellectual structure of a discipline, showing significant clustering of topically related authors
  - Shown to be effective in a broad range of disciplines, ranging from author cocitation analysis of scientific subfields to journal cocitation analysis of Economics and Marine Sciences



- Some researchers suggest that analyses of the full Web are no longer possible:
  - Such studies, attempting to characterize the “whole Web” are not feasible any more. ...the current estimated size of the indexable Web is around 2,200 million pages (*Judit Bar-Illan, Scientometrics 50(1) 2001*)
- This doesn't take into account use of new resources (such as Grids to harvest and process the data)

# Overview



- Web Characteristics
- Bibliometrics/Webometrics
- **Uses of Bibliometrics (link analysis)**
- New Approaches and New Uses

# Webometric analysis



- Link analysis
  - Use in search engines
- Co-citation analysis
  - Discovering web structures

# Search Engine Ranking: Link Analysis



- Assumptions:
  - If the pages pointing to this page are good, then this is also a good page
  - The words on the links pointing to this page are useful indicators of what this page is about
  - References: Page et al. 98, Kleinberg 98

# Ranking: Link Analysis



- Why does this work?
  - The official Toyota site will be linked to by lots of other official (or high-quality) sites
  - The best Toyota fan-club site probably also has many links pointing to it
  - Less high-quality sites do not have as many high-quality sites linking to them

# Ranking: PageRank



- Google uses the PageRank
- We assume page A has pages T1...Tn which point to it (i.e., are citations). The parameter d is a damping factor which can be set between 0 and 1. d is usually set to 0.85. C(A) is defined as the number of links going out of page A. The PageRank of a page A is given as follows:
  - $PR(A) = (1-d) + d (PR(T1)/C(T1) + \dots + PR(Tn)/C(Tn))$
  - Note that the PageRanks form a probability distribution over web pages, so the sum of all web pages' PageRanks will be one

# PageRank



- Similar to calculations used in scientific citation analysis (e.g., Garfield et al.) and social network analysis (e.g., Wasserman et al.)
- Similar to other work on ranking (e.g., the hubs and authorities of Kleinberg et al.)
- Computation is an iterative algorithm and converges to the principle eigenvector of the link matrix



# Early Work in Webometrics



- Larson (1996) used an adaptation of traditional co-citation analysis to examine the “intellectual structure of Cyberspace”
- This involved application of traditional co-citation analysis methods to the Web environment
- The assumption was that citation practice in the WWW environment served some of the same intellectual purposes as citation in the print world

# Traditional (Drexel-style) ACA



- Selection of authors,
- retrieval of co-citation frequencies,
- compilation of a raw co-citation matrix,
- conversion to a correlation matrix,
- multivariate analysis of correlation matrix (using principle components analysis, cluster analysis, and multidimensional scaling),
- interpretation: author's names on the maps can be translated into subject terms naming the clusters. Clusters and the dimensions (axes) on which they are placed reveal the intellectual structure of a field. *McCain (1990)*

# Example -- Geographic Sites



- Methodology
  - Search submitted to AltaVista (using the advanced search mode) was to find `link:pubweb.parc.xerox.com/map AND link:xtreme.gsfc.nasa.gov`, that is, to find a set of WWW documents containing links to both the Xerox Map browser, and the home page for NASA's AVHRR (Advanced Very High Resolution Radiometer) remote sensing projects.
  - Resulted (after throwing out irrelevant pages) in 43 sites that *co-cited* these sites

# Example Analysis



- The 43 sites contained 7209 outlinks to other sites, of these 332 existed in at least 3 of the 43 sites. These were examined and 34 “major” sites were selected as the core set for analysis

# Selected Sites



S1	The American Physical Society ( <a href="http://aps.org/">http://aps.org/</a> )
S2	NSF Geosciences Unidata Integrated Earth Information Server (IEIS) ( <a href="http://atm.geo.nsf.gov/">http://atm.geo.nsf.gov/</a> )
S3	Earth Observing System Home Page ( <a href="http://eos.nasa.gov/">http://eos.nasa.gov/</a> )
S4	WeatherNet ( <a href="http://cirrus.sprl.umich.edu/wxnet/">http://cirrus.sprl.umich.edu/wxnet/</a> )
S5	NASA-Goddard Climate and Radiation Branch ( <a href="http://climate.gsfc.nasa.gov/">http://climate.gsfc.nasa.gov/</a> )
S6	Center for Remote Sensing and Spatial Analysis ( <a href="http://deathstar.rutgers.edu/welcome.html">http://deathstar.rutgers.edu/welcome.html</a> )
S7	EcoWeb ( <a href="http://ecosys.drdr.virginia.edu/EcoWeb.html">http://ecosys.drdr.virginia.edu/EcoWeb.html</a> )
S8	EROS (Earth Resources Observation Systems) Data Center ( <a href="http://edcwww.cr.usgs.gov/eros-home.html">http://edcwww.cr.usgs.gov/eros-home.html</a> )
S9	The EnviroWeb- A Project of the EnviroLink Network ( <a href="http://envirolink.org/">http://envirolink.org/</a> )
S10	Live Access to Climate Data ( <a href="http://ferret.wrc.noaa.gov/ferret/main-menu.html">http://ferret.wrc.noaa.gov/ferret/main-menu.html</a> )
S11	Global Change Master Directory ( <a href="http://gcmd.gsfc.nasa.gov/">http://gcmd.gsfc.nasa.gov/</a> )
S12	Earth Science Division-Home Page ( <a href="http://geo.arc.nasa.gov/esd">http://geo.arc.nasa.gov/esd</a> )
S13	Information Center for the Environment (ICE), UC Davis ( <a href="http://ice.ucdavis.edu/">http://ice.ucdavis.edu/</a> )
S14	Climate Prediction Center ( <a href="http://nic.fb4.noaa.gov/">http://nic.fb4.noaa.gov/</a> )
S15	Xerox PARC Map Viewer ( <a href="http://pubweb.parc.xerox.com/map">http://pubweb.parc.xerox.com/map</a> )
S16	Public Use of Remote Sensing Data ( <a href="http://rsd.gsfc.nasa.gov/rsd/">http://rsd.gsfc.nasa.gov/rsd/</a> )
S17	SeaWiFS Project - Homepage ( <a href="http://seawifs.gsfc.nasa.gov/">http://seawifs.gsfc.nasa.gov/</a> )
S18	USGS Geo Data ( <a href="http://sun1.cr.usgs.gov/doc/edchome/ndcdb/ndcdb.html">http://sun1.cr.usgs.gov/doc/edchome/ndcdb/ndcdb.html</a> )
S19	Planet Earth Home Page ( <a href="http://www.nosc.mil/planet_earth/info.html">http://www.nosc.mil/planet_earth/info.html</a> )
S20	GeoWeb Home Page ( <a href="http://wings.buffalo.edu/geoweb/">http://wings.buffalo.edu/geoweb/</a> )
S21	Earthquake Information ( <a href="http://www.civeng.carleton.ca/cgi-bin/quakes">http://www.civeng.carleton.ca/cgi-bin/quakes</a> )
S22	ESRG (Earth Space Research Group) Gateway ( <a href="http://www.crseo.ucsb.edu/esrg.html">http://www.crseo.ucsb.edu/esrg.html</a> )
S23	EOS Volcanology ( <a href="http://www.geo.mtu.edu/eos/">http://www.geo.mtu.edu/eos/</a> )
S24	Mission to Planet Earth ( <a href="http://www.hq.nasa.gov/office/mtp/e/">http://www.hq.nasa.gov/office/mtp/e/</a> )
S25	Global Warming Update ( <a href="http://www.ncdc.noaa.gov/gblwrmupd/global.html">http://www.ncdc.noaa.gov/gblwrmupd/global.html</a> )
S26	The National Center for Geographic Information and Analysis ( <a href="http://www.ncgia.ucsb.edu/">http://www.ncgia.ucsb.edu/</a> )
S27	National Geophysical Data Center ( <a href="http://www.ngdc.noaa.gov/ngdc.html">http://www.ngdc.noaa.gov/ngdc.html</a> )
S28	NOAA National Oceanographic Data Center (NODC) Home Page ( <a href="http://www.nodc.noaa.gov/">http://www.nodc.noaa.gov/</a> )
S29	National Operational Hydrologic Remote Sensing Center ( <a href="http://www.nohrsc.nws.gov/">http://www.nohrsc.nws.gov/</a> )
S30	NOAA/PMEL/TAO El Nino Theme Page ( <a href="http://www.pmel.noaa.gov/to-ga-tao/el-nino/home.html">http://www.pmel.noaa.gov/to-ga-tao/el-nino/home.html</a> )
S31	NCAR Home Page ( <a href="http://www.ucar.edu/">http://www.ucar.edu/</a> )
S32	United States Geological Survey Home Page ( <a href="http://www.usgs.gov/">http://www.usgs.gov/</a> )
S33	Current Weather Maps/Movies ( <a href="http://wxweb.msu.edu/weather/">http://wxweb.msu.edu/weather/</a> )
S34	AVHRR Pathfinder Home Page ( <a href="http://xtreme.gsfc.nasa.gov/">http://xtreme.gsfc.nasa.gov/</a> )

# Example Analysis



- The 34 sites were analysed to construct a cocitation matrix
  - This stage requires the ability to search for “citing documents”, that is, those with links to the items in the core set and also the ability to conduct the many searches required (for any core set of size  $N$ , there are  $N^2/2 - N$  searches required -- one for each pair of items in the core set)
  - We developed a “robot” search program to submit these searches to AltaVista and collect the results
- The result was a cocitation matrix...

# Cocitation Matrix



Raw Cocitation Frequency	
ID	
31	
32	31
33	31 32
34	31 32 33
35	31 32 33 34
36	31 32 33 34 35
37	31 32 33 34 35 36
38	31 32 33 34 35 36 37
39	31 32 33 34 35 36 37 38
310	31 32 33 34 35 36 37 38 39
311	31 32 33 34 35 36 37 38 39 310
312	31 32 33 34 35 36 37 38 39 310 311
313	31 32 33 34 35 36 37 38 39 310 311 312
314	31 32 33 34 35 36 37 38 39 310 311 312 313
315	31 32 33 34 35 36 37 38 39 310 311 312 313 314
316	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315
317	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316
318	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317
319	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318
320	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318 319
321	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318 319 320
322	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318 319 320 321
323	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318 319 320 321 322
324	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318 319 320 321 322 323
325	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324
326	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325
327	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326
328	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327
329	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328
330	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329
331	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330
332	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330 331
333	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330 331 332
334	31 32 33 34 35 36 37 38 39 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330 331 332 333

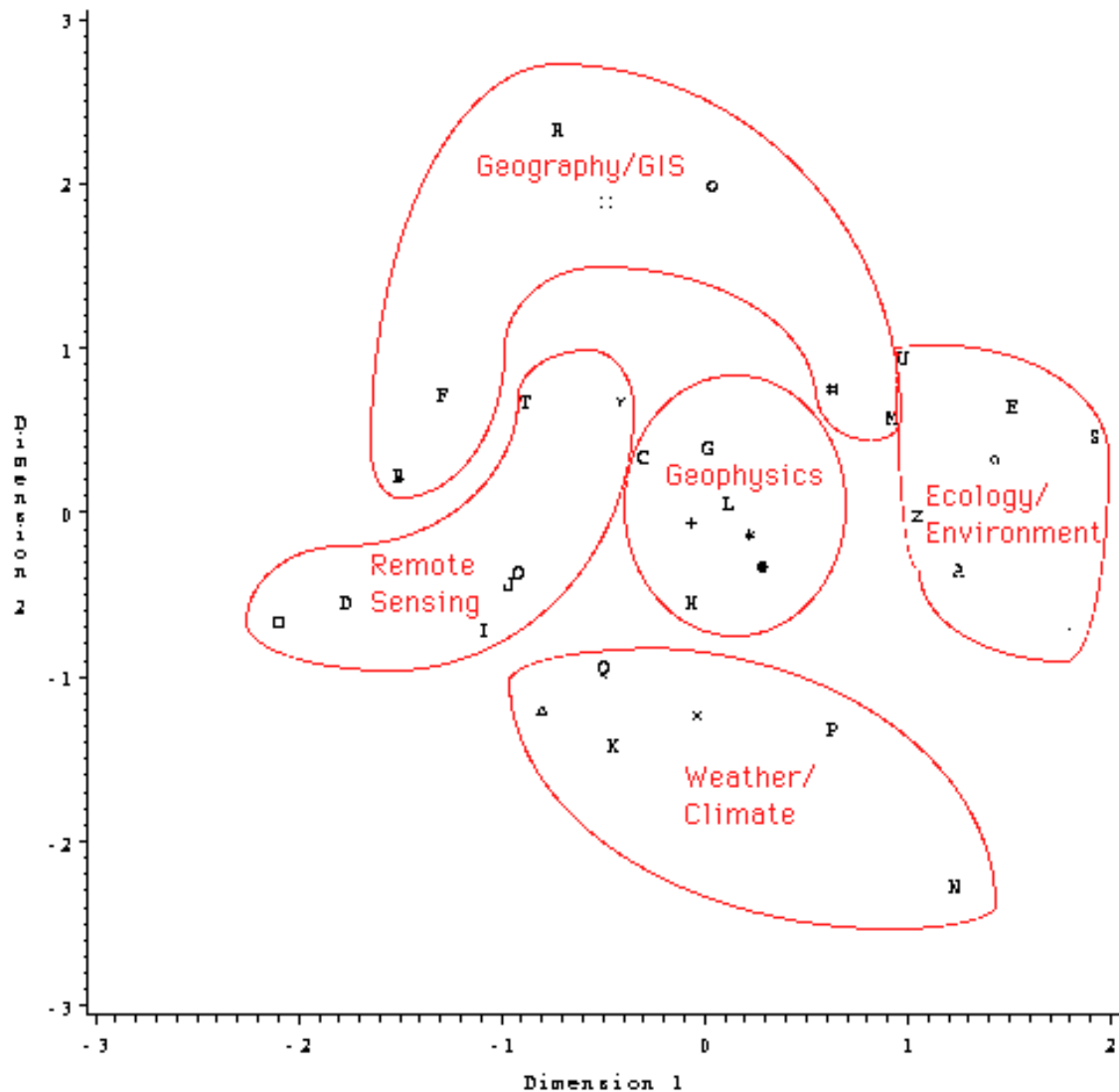
# Example Analysis



- The raw cocitation matrix was converted to a correlation matrix
- Next the correlation matrix was analysed using Multidimensional Scaling
  - MDS is a class of techniques for uncovering “hidden structure” in data.
- The result of applying MDS to a set of proximities results in the output of a spatial representation of the data consisting of a geometric configuration of points, in effect a mapping of the information onto a two or three dimensional plane.



# MDS Output



# Overview



- Web Characteristics
- Bibliometrics/Webometrics
- Uses of Bibliometrics (link analysis)
- **New Approaches and New Uses**

# Lin, White AuthorLink



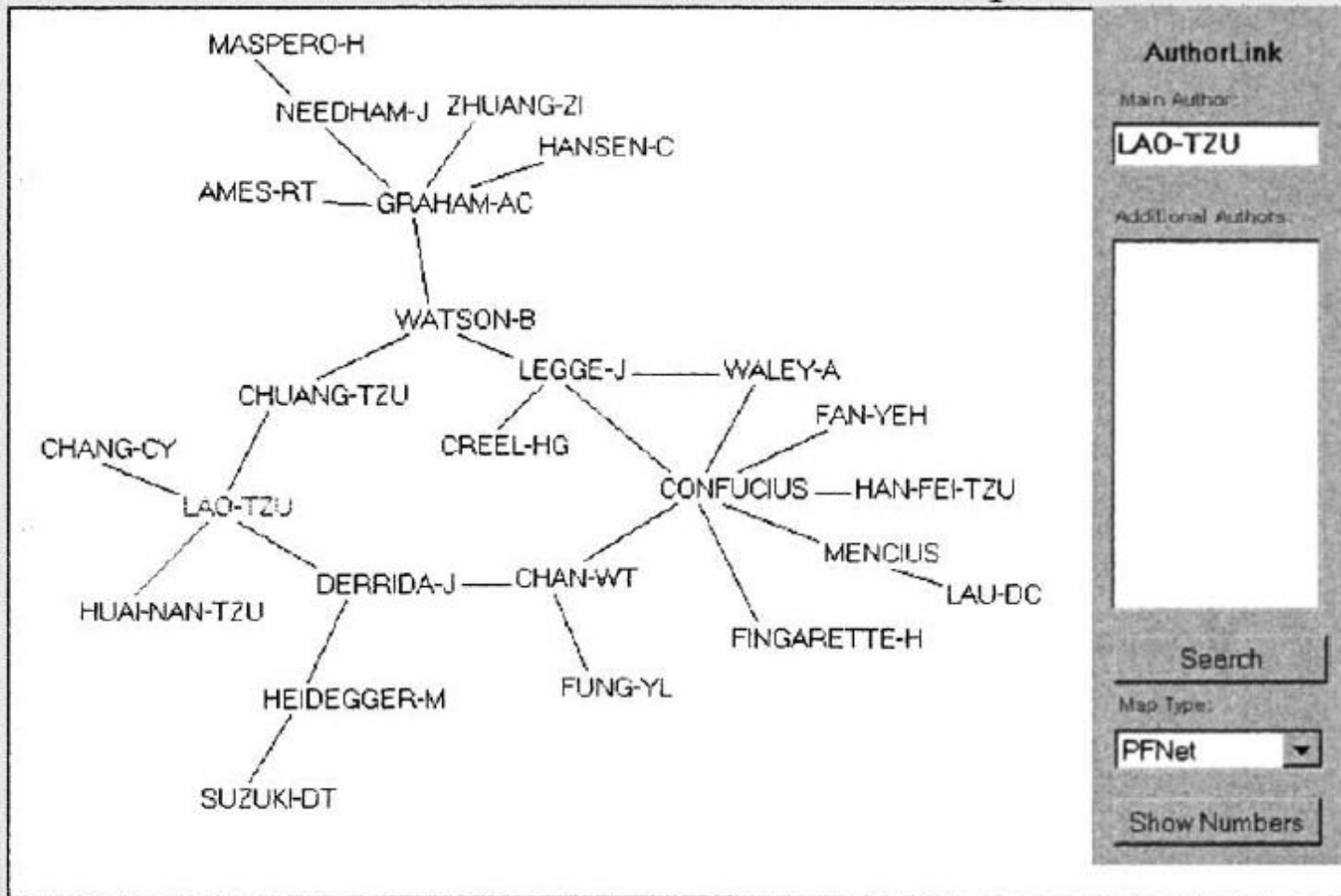
- Lin and White have developed a system for dynamic visualization of cocitation data called AuthorLink
- Provides support for analysis and visualization using the “traditional” tool, ISI’s Science Citation Index

# Lin, White AuthorLink Display



Author Search:

## Instant Author Co-citation Map

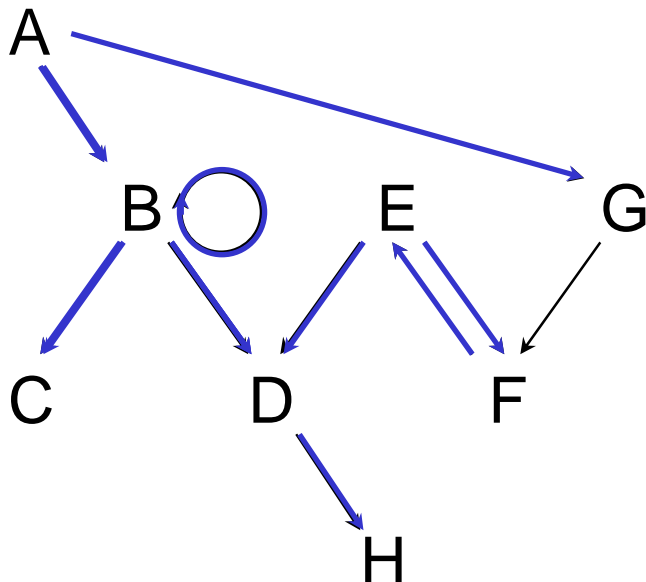


# Complex Link Analysis



- Ingwersen and his research group at the Royal Danish School of Library and Information Science are probably the central research group on Webometrics today. The following pages are derived from a presentation by Lennart Björneborn a Doctoral Student in this research group

# Basic Concepts



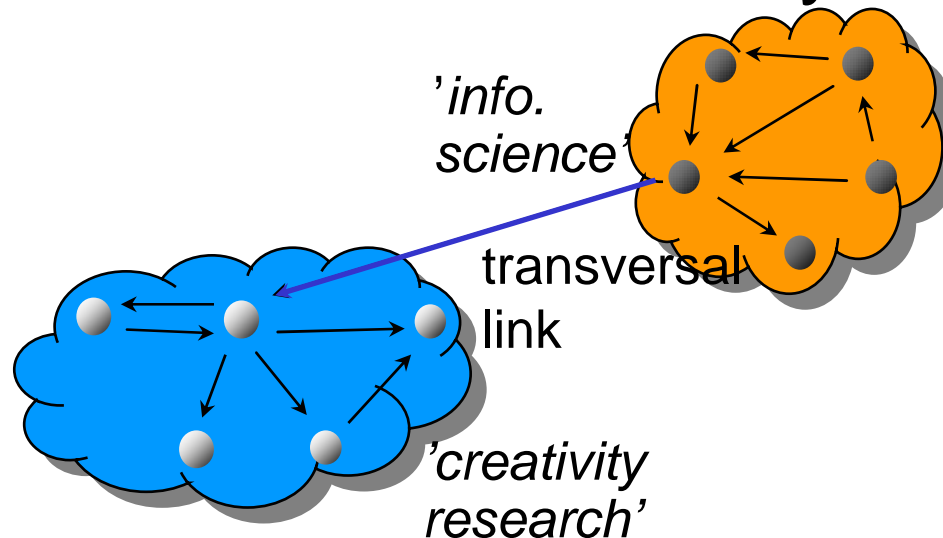
co-links

- B has an outlink to C : ~ reference
- B has an inlink from A : ~ citation
- B has a selflink : ~ self-citation
- E and F are reciprocally linked
- A is transitively linked with H via B - D
- A has a transversal link to G : short cut
- C and D are co-linked from B, i.e. shared inlinks: co-citation
- B and E are co-linking to D, i.e. shared outlinks: bibliog.coupling

# transversal links



- transversal links = short cuts or 'weak ties' between otherwise 'distant' web clusters (e.g., subject domains, interest communities)
- example 1: info.sci. researcher with transversal link to creativity research





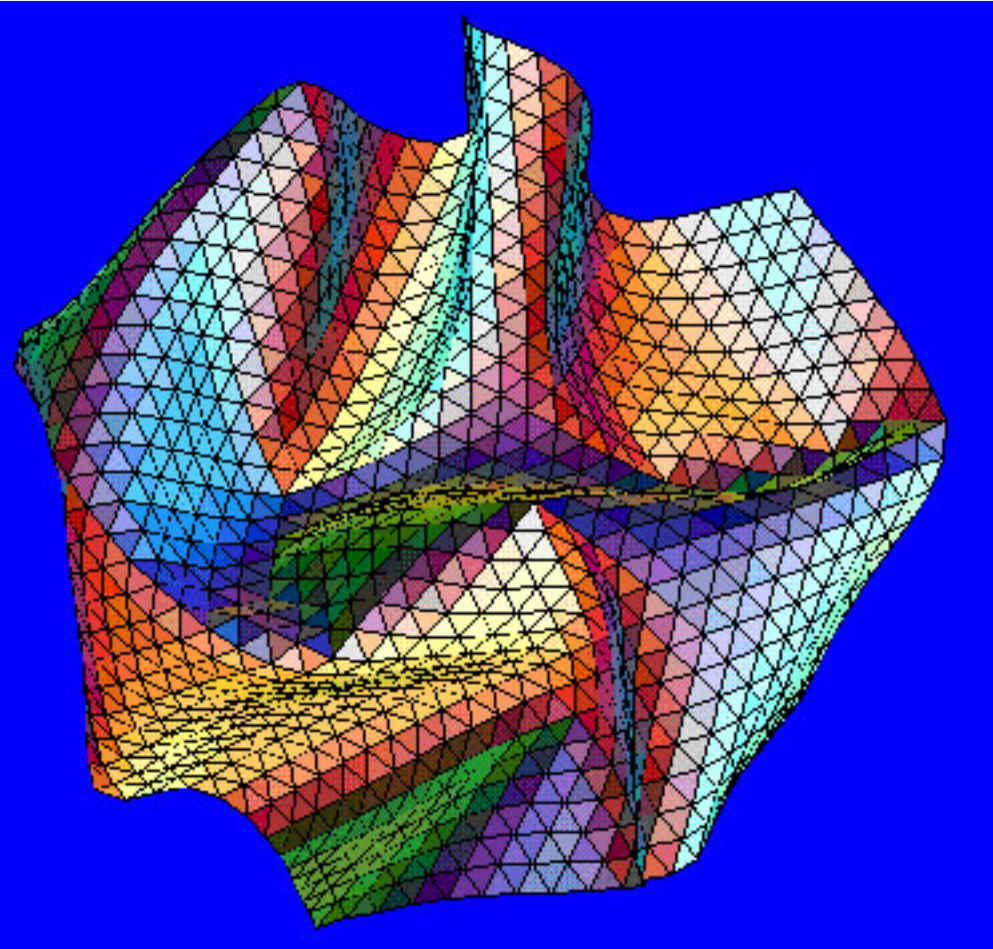
- transversal links may create small-world phenomena
  - in the shape of short distances between nodes in the Web graph
  - making the Web more strongly connected and ‘crumpled up’



# 'crumpled-up' Web

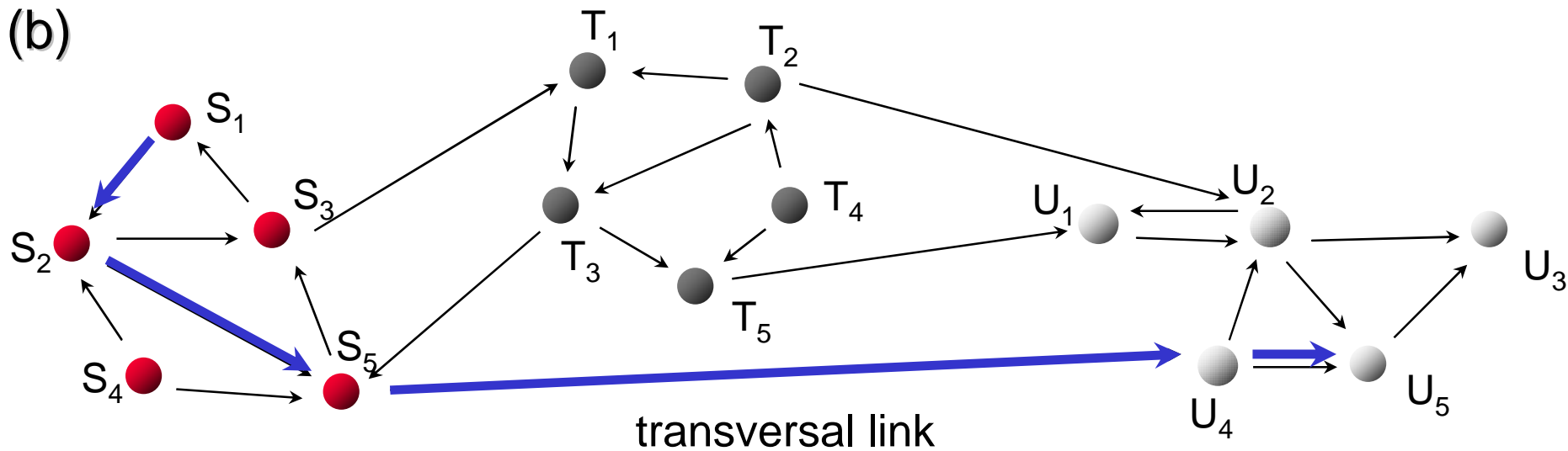
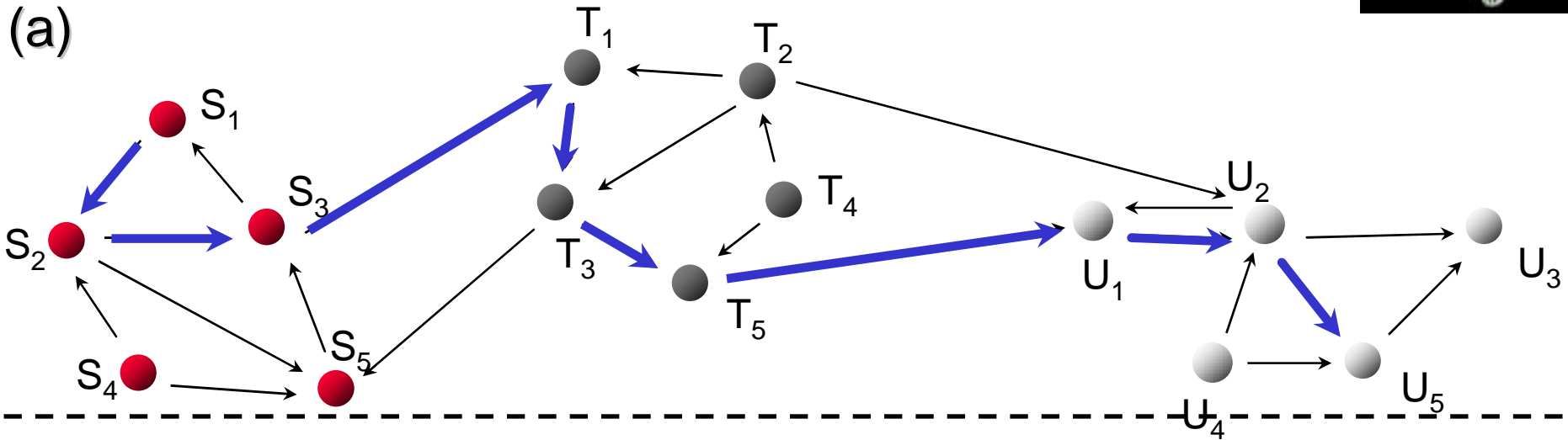


Image: <http://jfi.uchicago.edu/~tten/rainbow/Crumpling/>



- analogy: crumpled-up paper
- paths between selected points get shorter as opposite parts of structure are pulled together
- paper = three-dimensional
- www = billion-dimensional
- every new web outlink may interconnect with any existing and accessible web node
  - creating a new dimension in the 'crumpled-up' Web space
- every new link = 'hook'
  - stretching and reshaping existing web network

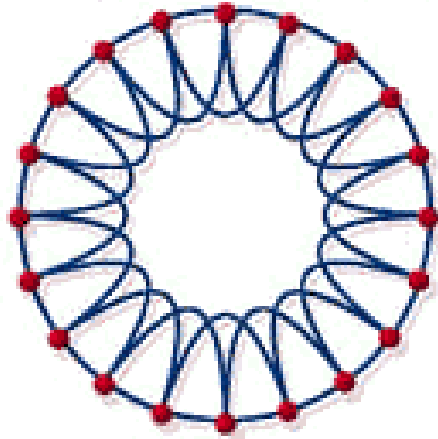
# Traversal Links



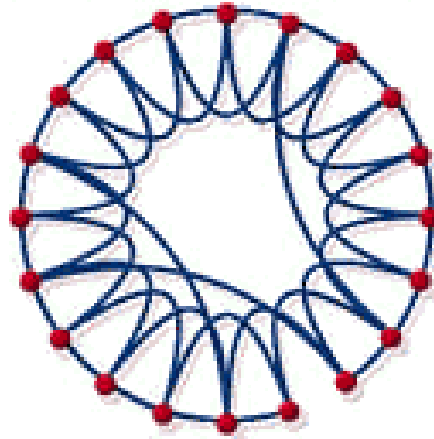
# small-world networks



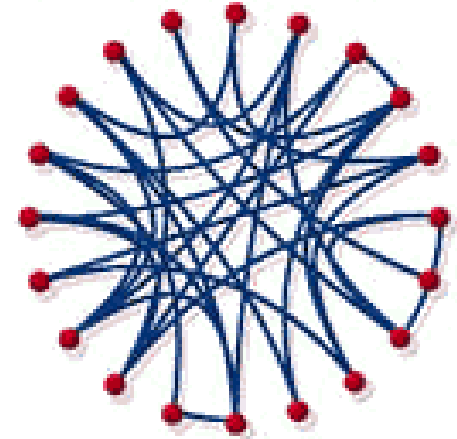
**Regular network:**  
connections to 4 nearest neighbours



**Small-world network:**  
a few long-range connections



**Random network:**  
points connected haphazardly

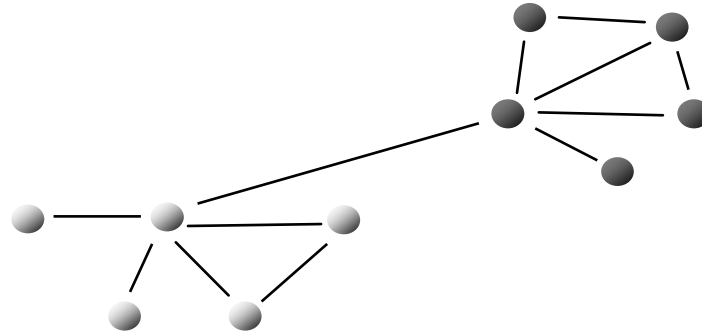


Increasing randomness

Fig.: Matthews (1998), *New Scientist*

- small-world = highly clustered + short paths [graph theory]
- small-world = order + diversity = complementarities
- short distances through short cuts between nodes in network
- small-world = short local + short global distances

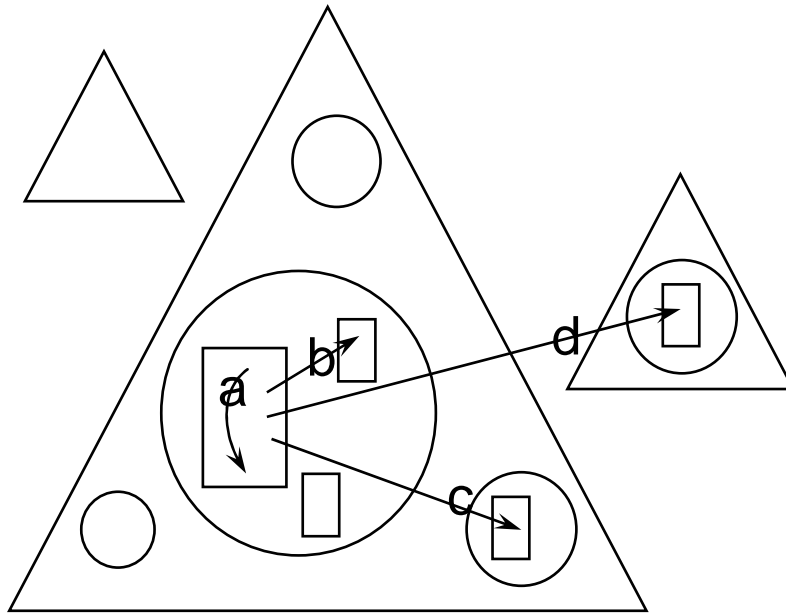
# Small-World phenomena


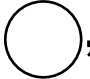



- small-world theories in social network analysis in 1960s
  - small-world phenomena in the shape of short distances between two arbitrary persons through intermediate chains of acquaintances
- ‘six degrees of separation’ (cf. play 1990, movie 1993)
  - cf. ‘Six degrees of Kevin Bacon’ (<http://www.cs.virginia.edu/oracle>)
- ‘small worlds’ in biological, technological and social networks: neural networks, electrical power grids, diffusion of ideas and epidemics,...
- lack of research on small-world phenomena in informational networks: WWW, citation databases, semantic networks, thesauri, ...



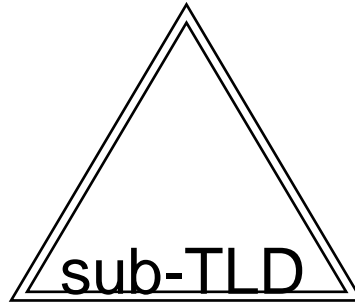
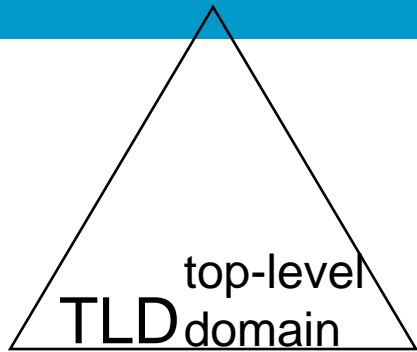
# basic levels of web nodes



- 3 basic levels of web nodes: pages , domains , TLDs 
- different levels of selflinks and outlinks
  - a = page selflink
  - b = page outlink and domain selflink
  - c = domain outlink and TLD selflink
  - d = TLD outlink
- more levels: frames, directories, sub-domains, sub-TLDs ...

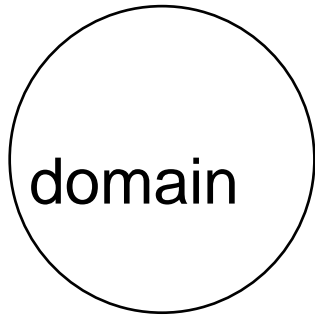
top level domains

# more web node notation

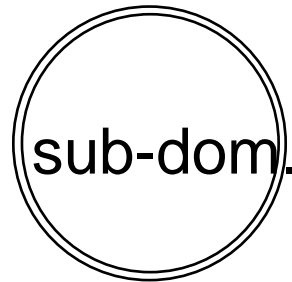


Example: .uk

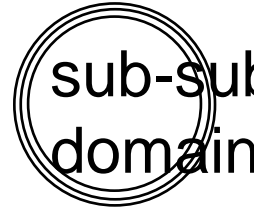
.ac.uk



.wlv.ac.uk

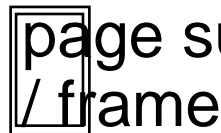
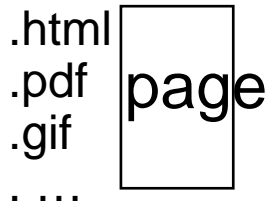


.scit.wlv.ac.uk



x1.scit.wlv.ac.uk

sub-sub-sub-...

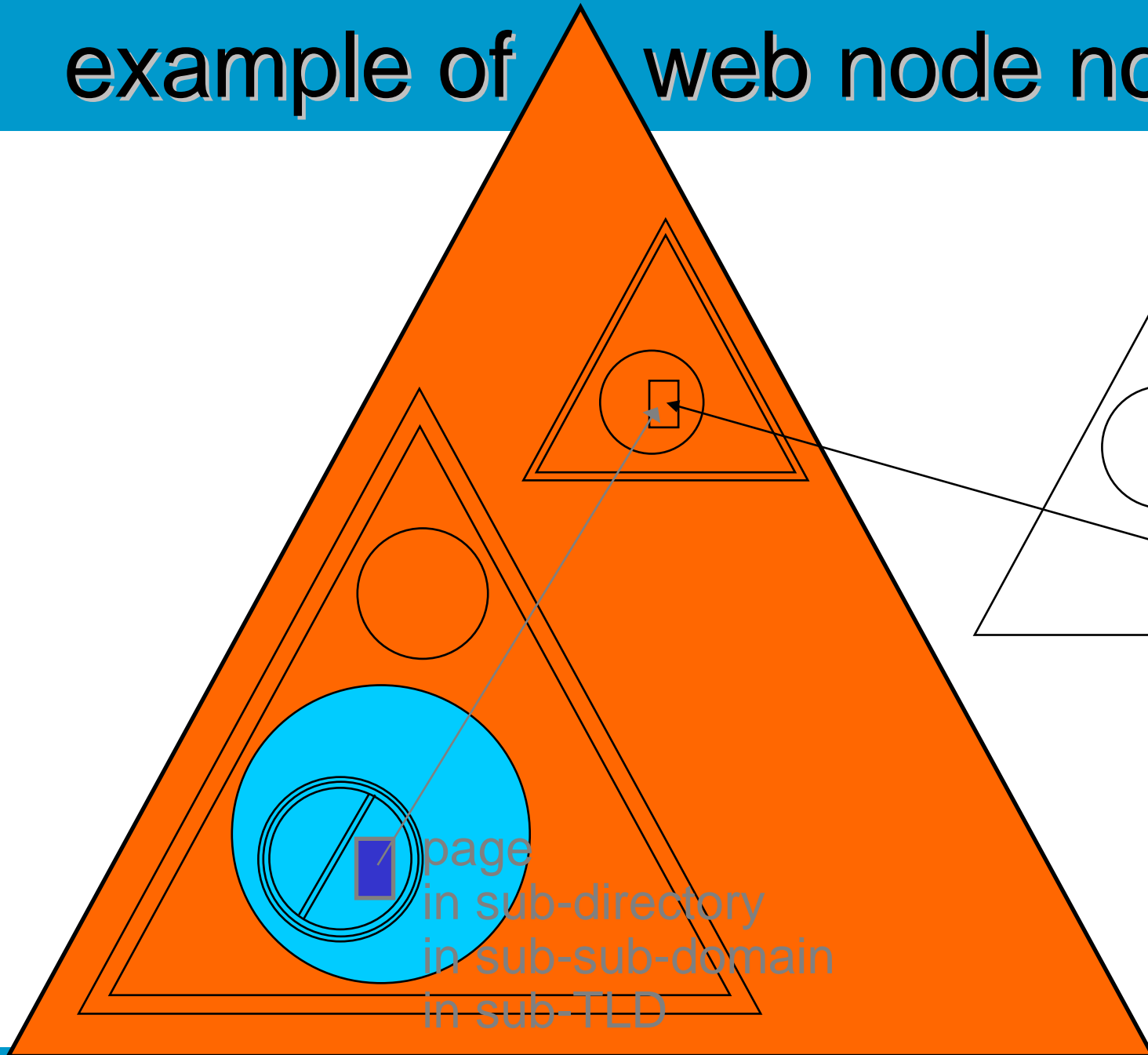


x1.scit.wlv.ac.uk/yy



.scit.wlv.ac.uk/zz/aa/bb

# example of web node notation



page  
in sub-directory  
in sub-sub-domain  
in sub-TLD

.edu

.co.uk

sub-domain: .scit.wlv.ac.uk

domain: .wlv.ac.uk

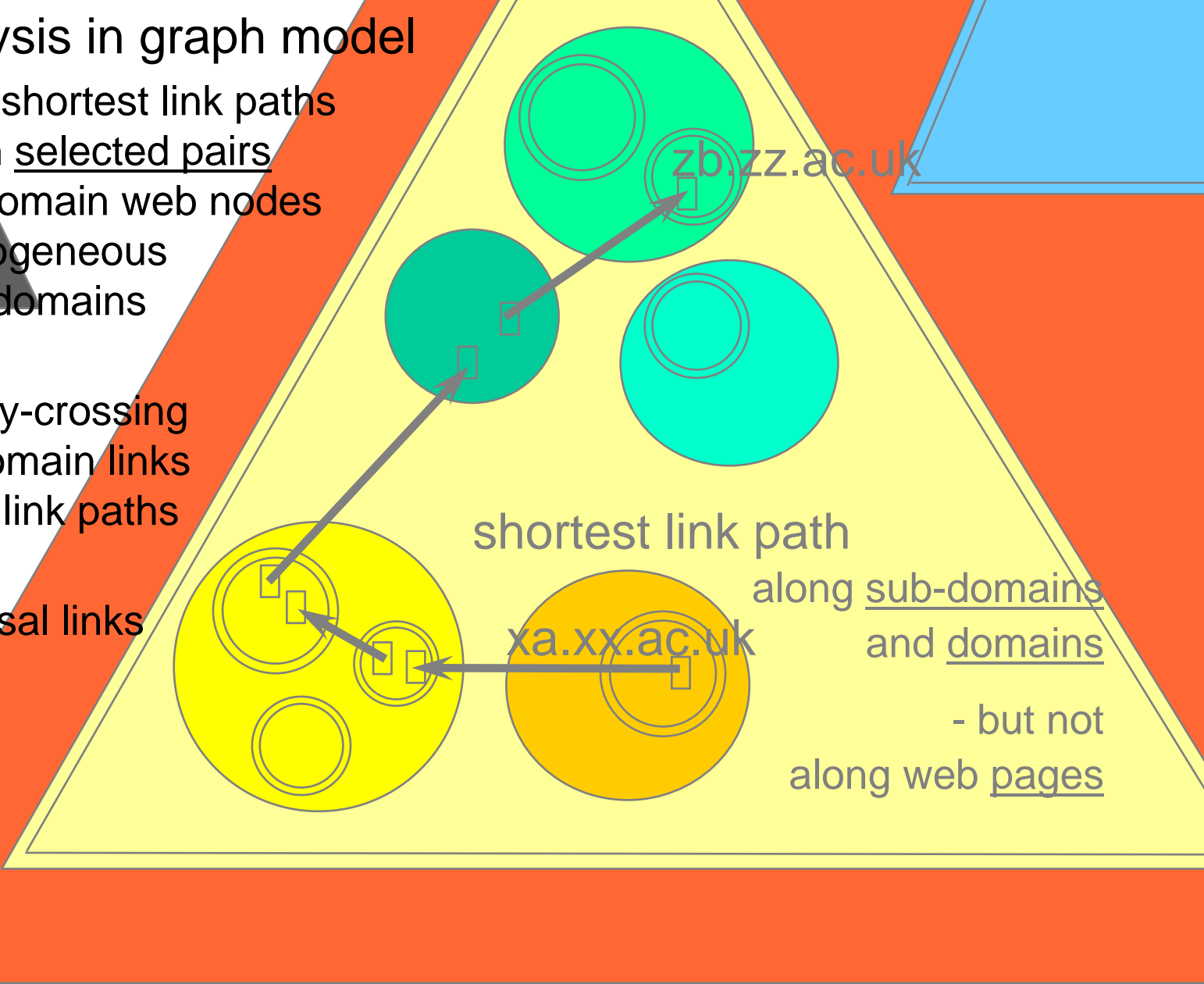
sub-TLD: .ac.uk

TLD: .uk

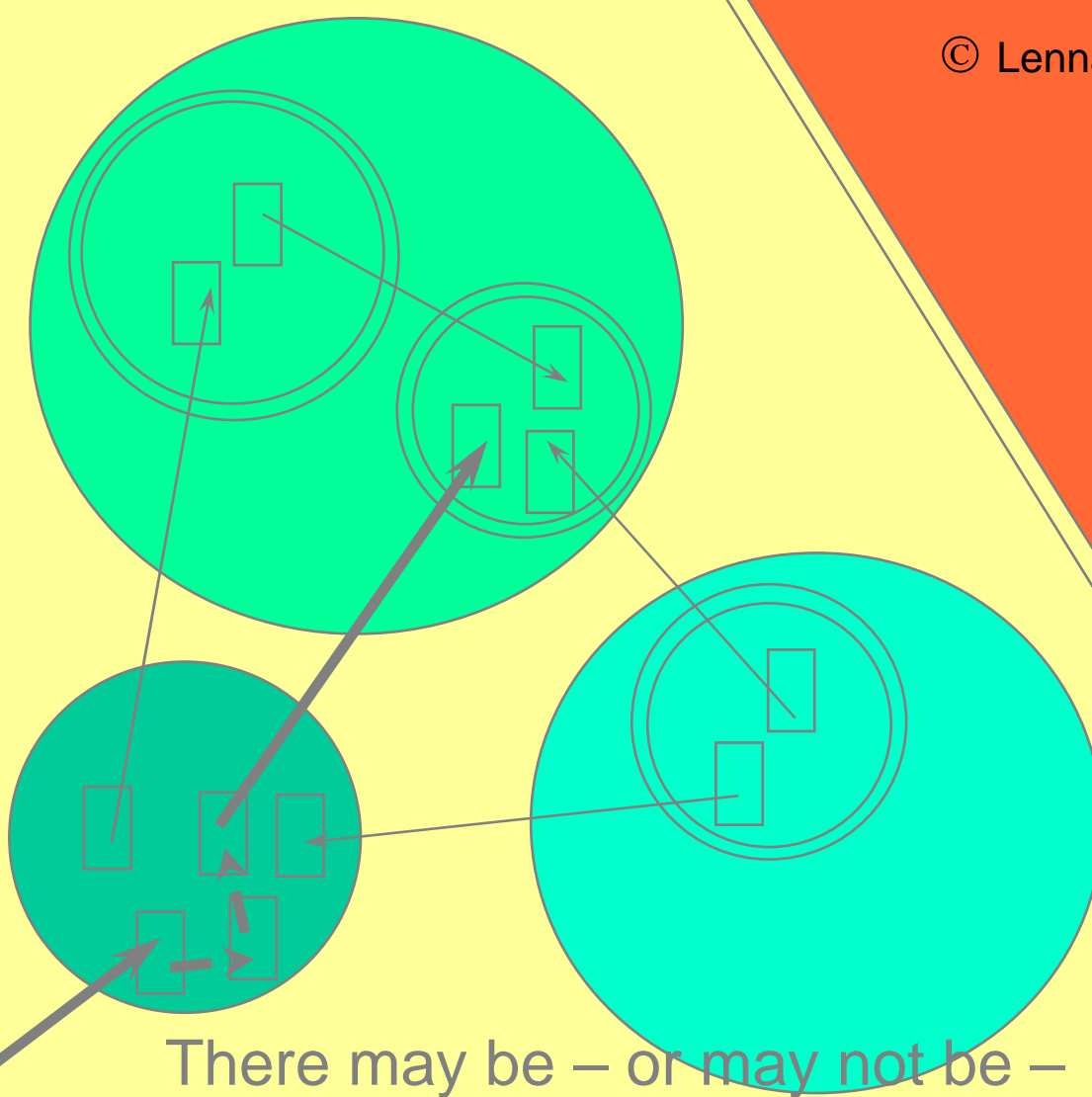


# path analysis in graph model

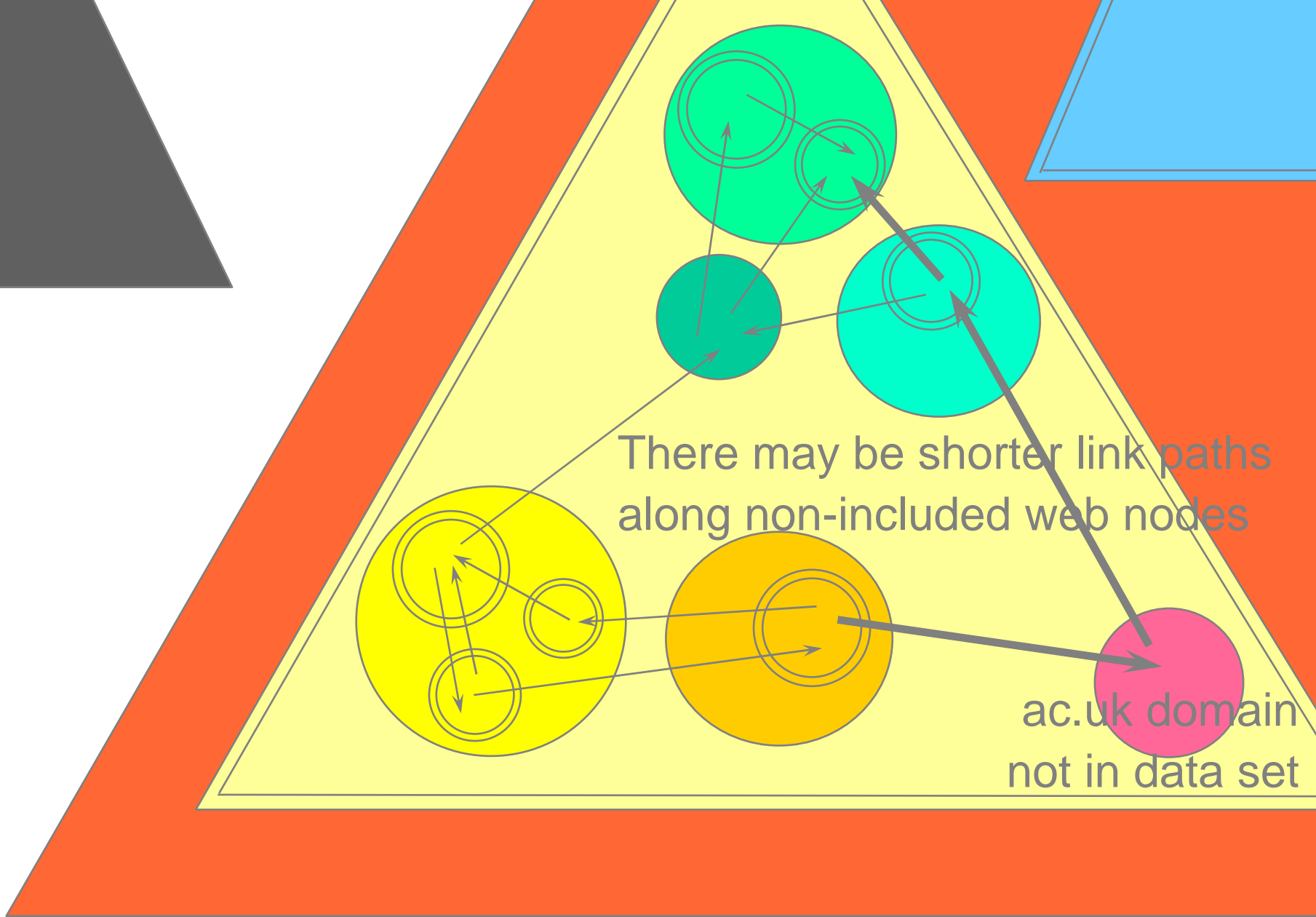
- analyse shortest link paths between selected pairs of sub-domain web nodes in heterogeneous subject domains
- identify boundary-crossing cross-domain links on such link paths
- possible transversal links

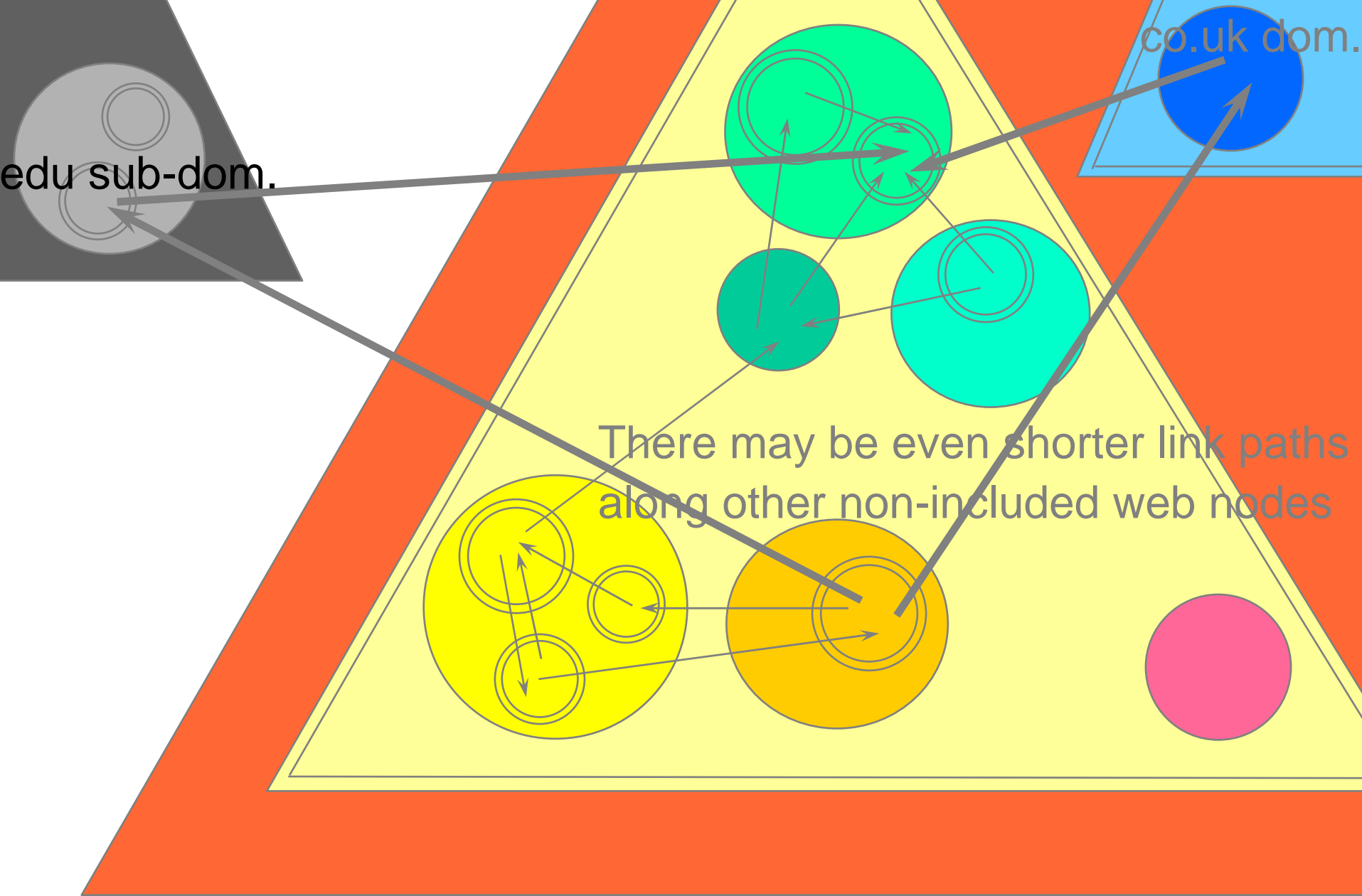


shortest link path  
along sub-domains  
and domains  
- but not  
along web pages



There may be – or may not be –  
local link paths connecting web pages  
on the overall link path  
along web domains and sub-domains



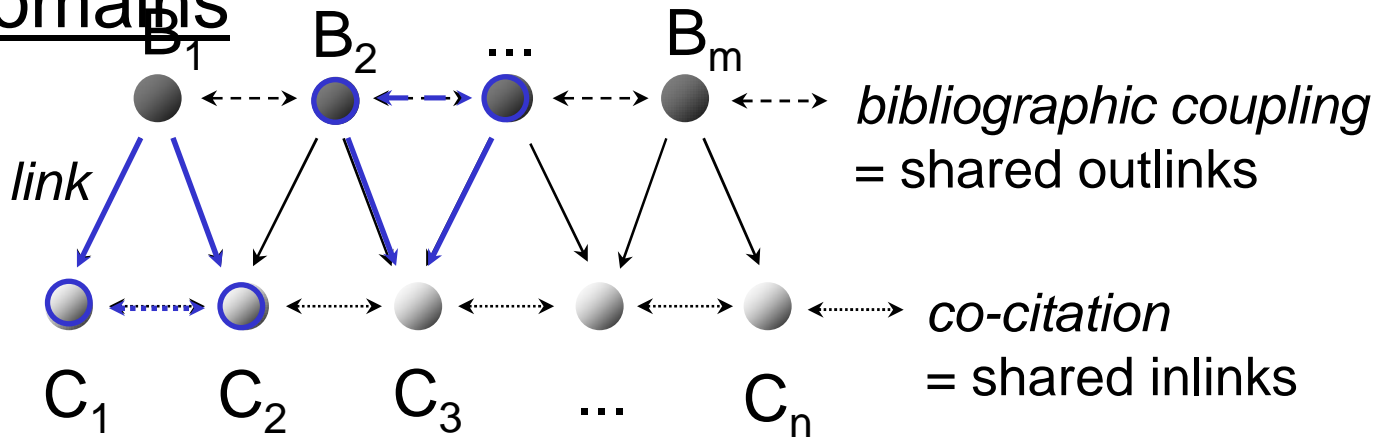


# shared outlinks in link data set



Björneborn 2001-2002

- analysing shared outlinks between domains / sub-  
domains



- shared outlinks between web pages with bookmark lists
  - researchers' diverse interests, preferences and actions on the Web
  - may reflect emerging 'research fronts' or 'invisible colleges'
  - correlation with co-citations/bib.couplings in Science Citation Index?
- web mining 'undiscovered public knowledge' (cf. Swanson 1986)
- possible power law in distribution of shared outlinks

# co-linkage chain

- case studies of transversal **co-links**



researchers' homepages

researchers' bookmark lists

Jon Kleinberg: link structures,  
small-world phenomena

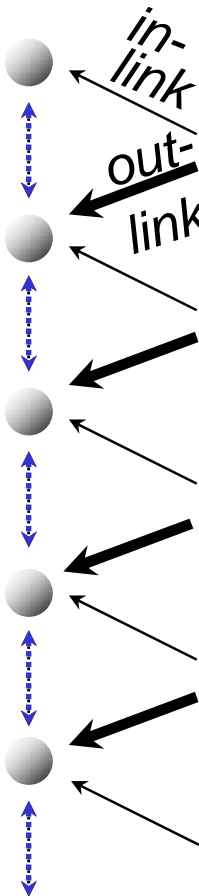
*co-citation*

Frances Heylighen: cybernetics

David Chalmers: philosophy of mind

Mark Warschauer: language/literacy

Michael Barlow: linguistics



Cliff Joslyn:

distributed knowledge systems

*bibliographic coupling*

Liane Gabora: interdiscipl. studies

Don Sorsa: education research

James Wong: computer-assisted  
language learning

NN

AltaVista's Advanced Search:  
(title:bookmark\* OR title:link\*) AND  
anchor:"Michael Barlow"

# users' information behavior



Björneborn 2001-2002

## traditional

### LIS research areas

- 'convergent' behavior
- goal-directed, rational
- e.g. Boolean searching
- present, explicit info.needs
- problems, work tasks
- 'information recovery'
- order = topic-focused, 'convergent' link structures

## complementary

### LIS research areas

- 'divergent' behavior
- non-goal-directed, intuitive
- e.g. browsing, serendipity
- latent, implicit info.needs
- triggered interests, curiosity
- 'information discovery'
- diversity = topic-scattered, 'divergent' link structures

- *which* combinations of link structures may facilitate *which* combinations of behavior?

# exploratory capabilities



- LB proposal: broadened aim of LIS research:
- "facilitating users to explore and exploit options embedded in information environments"
- small-world link structures may enhance exploratory capabilities in an informational network such as the Web
  - including chances for serendipity - encountering the unexpected
- easier to traverse and explore an information space if both local and global distances are short = 'small world'
- small-world link structures may facilitate combinations of 'convergent' and 'divergent' information behavior



# possible implications



- webometric studies of small-world link structures may provide a better understanding of complex topologies, functionalities and potentials of the web
- may be utilised in
  - social informatics – exploring emerging cultural/social formations
  - digital libraries – making them as serendipity-prone as physical libraries...
    - facilitating exploratory capabilities!
  - creativity research – networked knowledge creation and diffusion
  - web mining – identifying fertile areas for cross-disciplinary exploration
  - search engines – more exhaustive web traversal + harvesting + ranking
  - browsers – visualisation / navigation facilities stimulating serendipity
  - etc.

# Future



- With massively parallel systems (Grid computing) and very-large scale storage (the Data Grid) it *should* be possible to do very-large-scale analyses of the structure and content of the WWW
- Perhaps we finally begin to see the “intellectual structure of Cyberspace”