

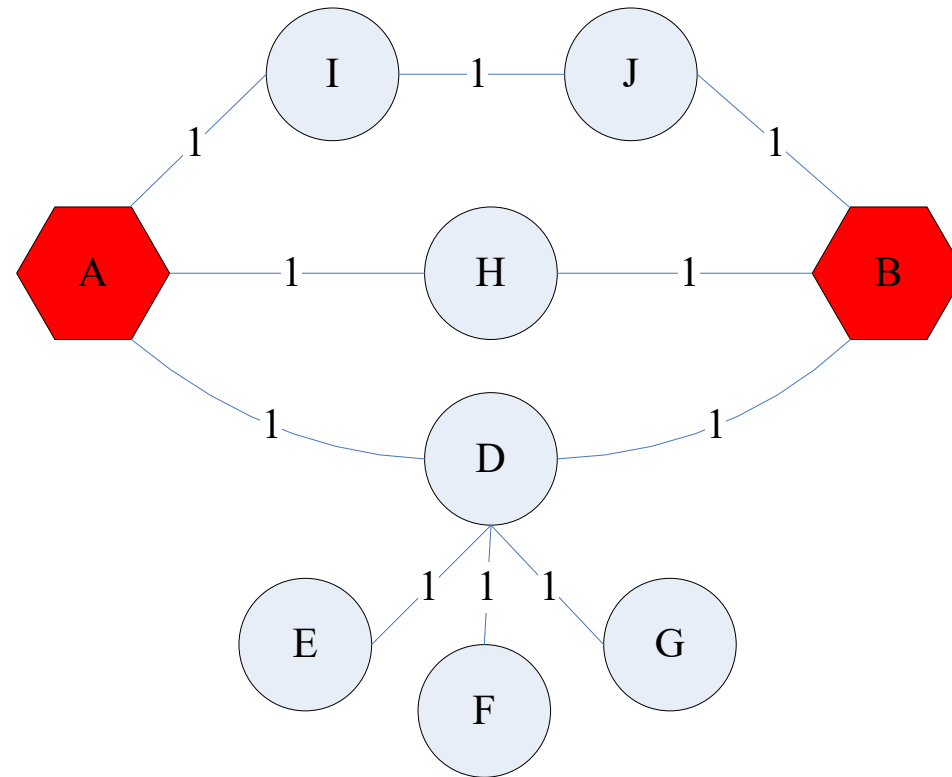
# Discovering the Topic Vector S

- **Create different PageRanks for different topics**
  - The 16 DMOZ top-level categories:
    - arts, business, sports,...
- **Which topic ranking to use?**
  - User can pick from a menu
  - Classify query into a topic
  - Can use the **context** of the query
    - E.g., query is launched from a web page talking about a known topic
    - History of queries e.g., “basketball” followed by “Jordan”
  - User context, e.g., user’s bookmarks, ...

# Application to Measuring Proximity in Graphs

Random Walk with Restarts:  $S$  is a single element

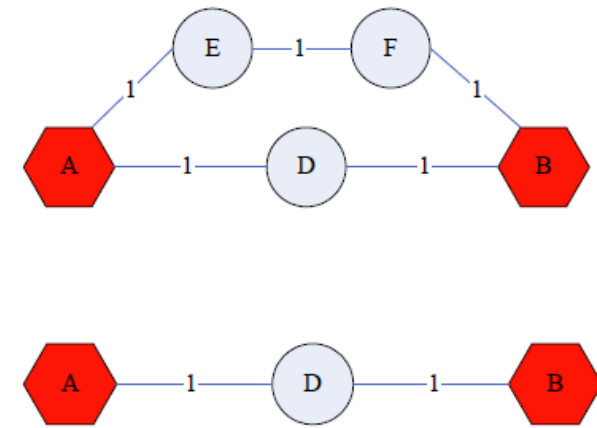
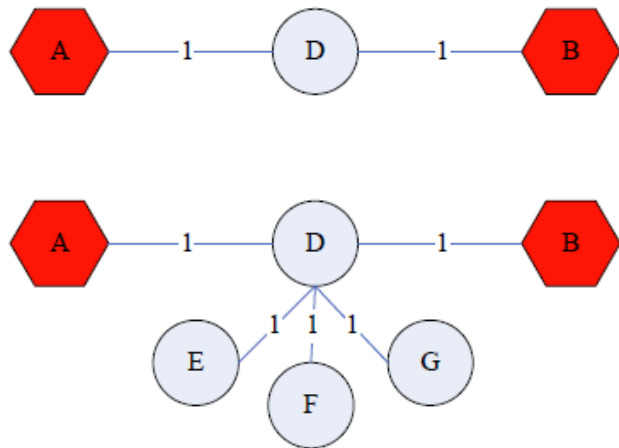
# Proximity on Graphs



**a.k.a.: Relevance, Closeness, 'Similarity'...**

# Good proximity measure?

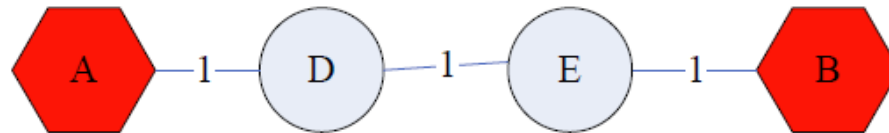
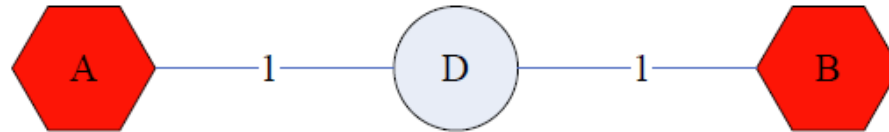
- Shortest path is not good:



- No effect of degree-1 nodes (E, F, G)!
- Multi-faceted relationships

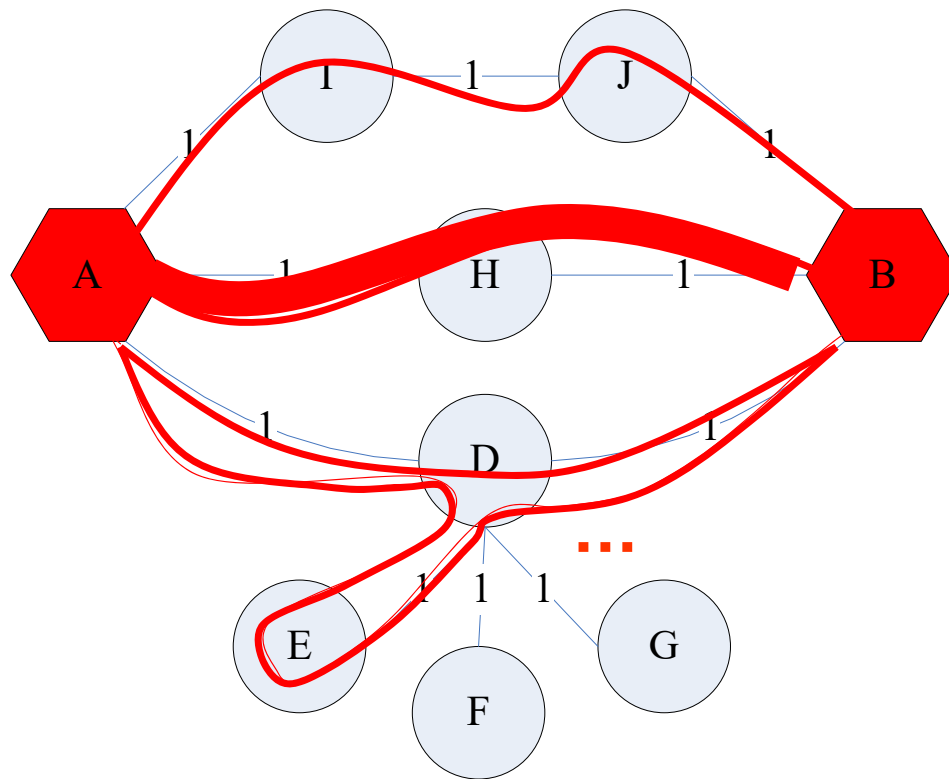
# Good proximity measure?

- Network flow is not good:



- Does not punish long paths

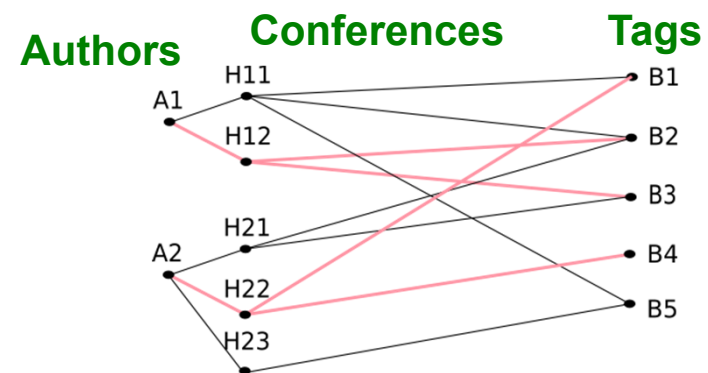
# What is good notion of proximity?



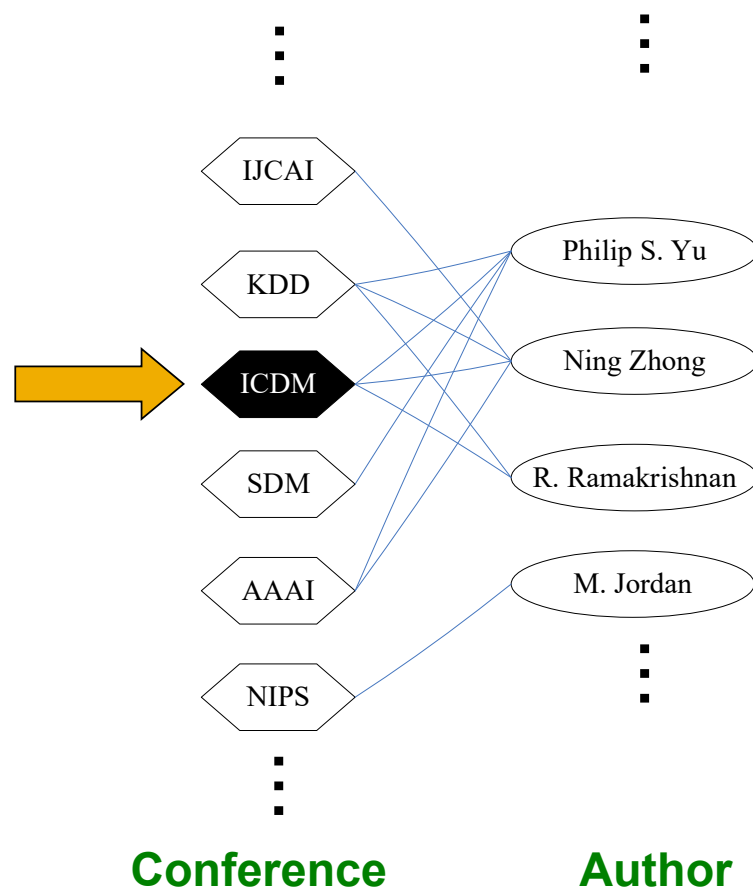
- Multiple connections
- Quality of connection
  - Direct & Indirect connections
  - Length, Degree, Weight...

# SimRank: Idea

- **SimRank:** Random walks from a **fixed node** on  **$k$ -partite graphs**
- **Setting:**  **$k$ -partite graph** with  **$k$  types of nodes**
  - E.g.: Authors, Conferences, Tags
- **Topic Specific PageRank** from node  **$u$** : **teleport set  $S = \{u\}$**
- Resulting scores measures similarity to node  **$u$**
- **Problem:**
  - Must be done once for each node  **$u$**
  - Suitable for sub-Web-scale applications



# SimRank: Example

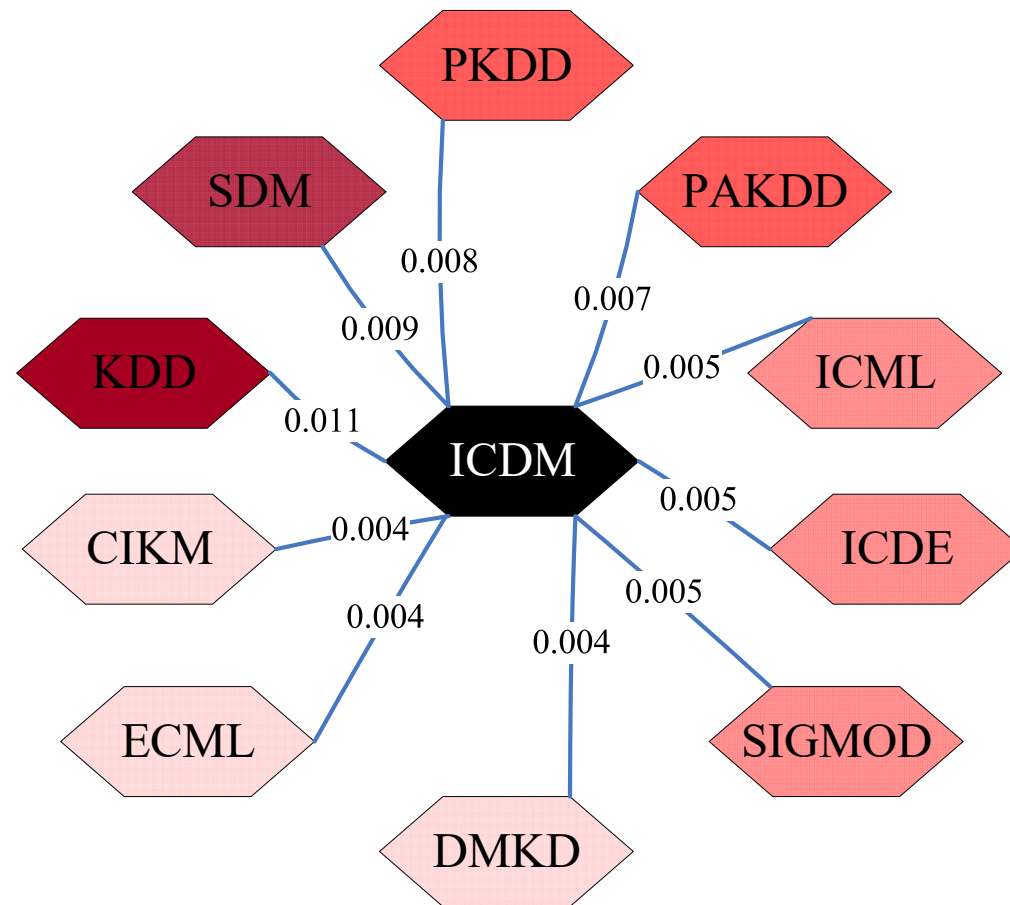


**Q:** What is most related conference to **ICDM**?

**A:** Topic-Specific PageRank with teleport set  $S=\{\text{ICDM}\}$



# SimRank: Example



# PageRank: Summary

## ■ “Normal” PageRank:

- Teleports uniformly at random to any node
- All nodes have the same probability of surfer landing there:  $\mathbf{S} = [0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1, 0.1]$

## ■ Topic-Specific PageRank also known as Personalized PageRank:

- Teleports to a topic specific set of pages
- Nodes can have different probabilities of surfer landing there:  $\mathbf{S} = [0.1, 0, 0, 0.2, 0, 0, 0.5, 0, 0, 0.2]$

## ■ Random Walk with Restarts:

- Topic-Specific PageRank where teleport is always to the same node.  $\mathbf{S} = [0, 0, 0, 0, \mathbf{1}, 0, 0, 0, 0, 0]$

# TrustRank: Combating the Web Spam

# What is Web Spam?

- **Spamming:**

- Any deliberate action to boost a web page's position in search engine results, incommensurate with page's real value

- **Spam:**

- Web pages that are the result of spamming

- This is a very broad definition

- Search engine optimization (SEO) industry might disagree!

- Approximately **10-15%** of web pages are spam

# Web Search

## ■ Early search engines:

- Crawl the Web
- Index pages by the words they contained
- Respond to search queries (lists of words) with the pages containing those words

## ■ Early page ranking:

- Attempt to order pages matching a search query by “importance”
- **First search engines considered:**
  - (1) Number of times query words appeared
  - (2) Prominence of word position, e.g. title, header

# First Spammers

- As people began to use search engines to find things on the Web, those with commercial interests tried to **exploit search engines** to bring people to their own site – whether they wanted to be there or not
- **Example:**
  - Shirt-seller might pretend to be about “movies”
- **Techniques for achieving high relevance/importance for a web page**

# First Spammers: Term Spam

- **How do you make your page appear to be about movies?**
  - (1) Add the word movie 1,000 times to your page
  - Set text color to the background color, so only search engines would see it
  - (2) Or, run the query “movie” on your target search engine
  - See what page came first in the listings
  - Copy it into your page, make it “invisible”
- **These and similar techniques are term spam**

# Google's Solution to Term Spam

- Believe what people say about you, rather than what you say about yourself
  - Use words in the anchor text (words that appear underlined to represent the link) and its surrounding text
- PageRank as a tool to measure the “importance” of Web pages



# Why It Works?

## ■ Our hypothetical shirt-seller looses

- Saying he is about movies doesn't help, because others don't say he is about movies
- His page isn't very important, so it won't be ranked high for shirts or movies

## ■ Example:

- Shirt-seller creates 1,000 pages, each links to his with “movie” in the anchor text
- These pages have no links in, so they get little PageRank
- So the shirt-seller can't beat truly important movie pages, like IMDB

# Why it does not work?

Google™ [Web](#) [Images](#) [Groups](#) [News](#) [Froogle](#) [Local](#) [more »](#)

[Advanced Search](#)  
[Preferences](#)

---

**Web** Results 1 - 10 of about 969,000 for [miserable failure](#). (0.06 seconds)

[Biography of President George W. Bush](#)  
Biography of the president from the official White House web site.  
[www.whitehouse.gov/president/gwbbio.html](#) - 29k - [Cached](#) - [Similar pages](#)  
[Past Presidents](#) - [Kids Only](#) - [Current News](#) - [President](#)  
[More results from www.whitehouse.gov »](#)

[Welcome to MichaelMoore.com!](#)  
Official site of the gadfly of corporations, creator of the film Roger and Me and the television show The Awful Truth. Includes mailing list, message board, ...  
[www.michaelmoore.com/](#) - 35k - [Sep 1, 2005](#) - [Cached](#) - [Similar pages](#)

[BBC NEWS | Americas | 'Miserable failure' links to Bush](#)  
Web users manipulate a popular search engine so an unflattering description leads to the president's page.  
[news.bbc.co.uk/2/hi/americas/3298443.stm](#) - 31k - [Cached](#) - [Similar pages](#)

[Google's \(and Inktomi's\) Miserable Failure](#)  
A search for **miserable failure** on Google brings up the official George W. Bush biography from the US White House web site. Dismissed by Google as not a ...  
[searchenginewatch.com/sereport/article.php/3296101](#) - 45k - [Sep 1, 2005](#) - [Cached](#) - [Similar pages](#)



# Google vs. Spammers: Round 2!

- Once Google became the dominant search engine, spammers began to work out ways to fool Google
- **Spam farms** were developed to concentrate PageRank on a single page
- **Link spam:**
  - Creating link structures that boost PageRank of a particular page



# Link Spamming

- **Three kinds of web pages from a spammer's point of view**
  - **Inaccessible pages**
  - **Accessible pages**
    - e.g., blog comments pages
    - spammer can post links to his pages
  - **Owned pages**
    - Completely controlled by spammer
    - May span multiple domain names

# Link Farms

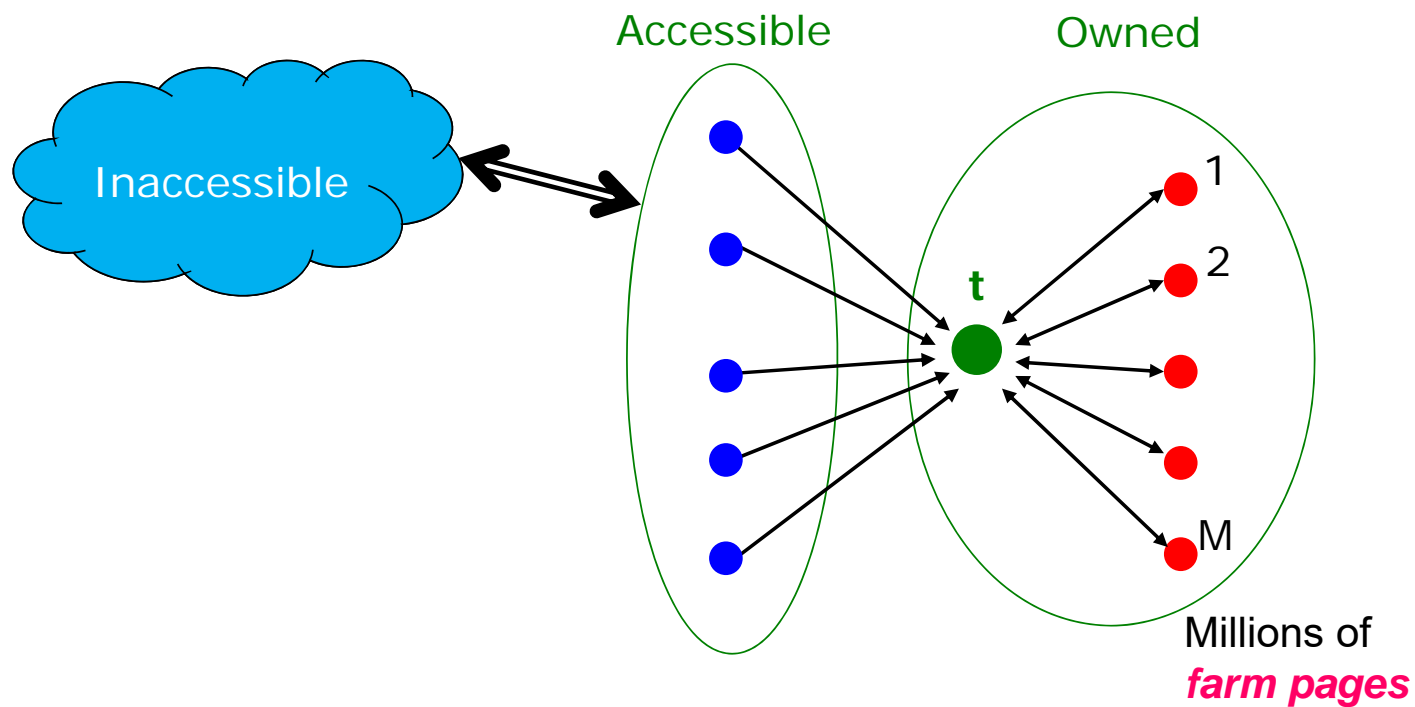
- **Spammer's goal:**

- Maximize the PageRank of target page  $t$

- **Technique:**

- Get as many links from accessible pages as possible to target page  $t$
- Construct “link farm” to get PageRank multiplier effect

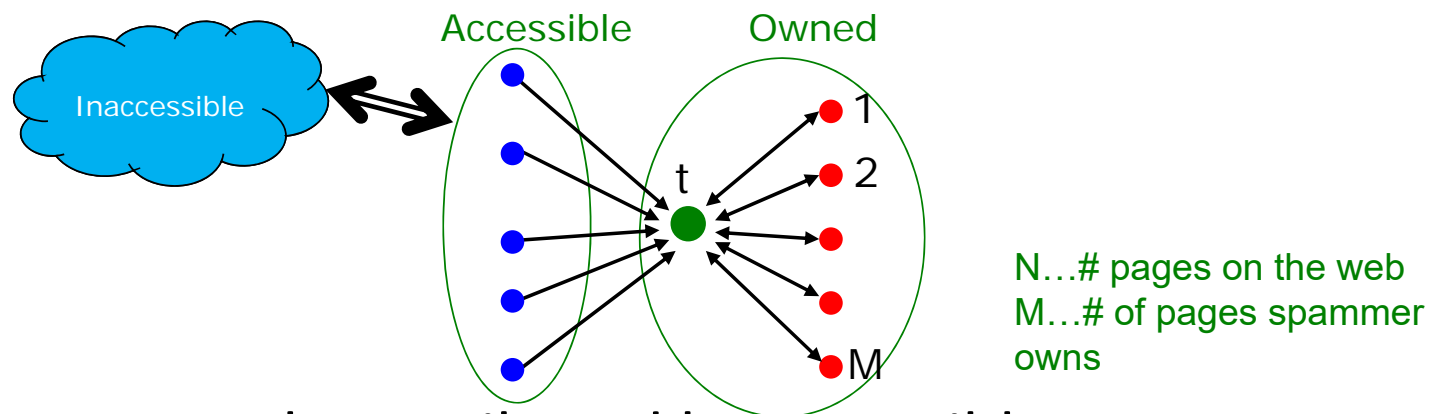
# Link Farms



One of the most common and effective organizations for a link farm



# Analysis



- $x$ : PageRank contributed by accessible pages

- $y$ : PageRank of target page  $t$

- Rank of each “farm” page =  $\frac{\beta y}{M} + \frac{1-\beta}{N}$

- $y = x + \beta M \left[ \frac{\beta y}{M} + \frac{1-\beta}{N} \right] + \frac{1-\beta}{N}$

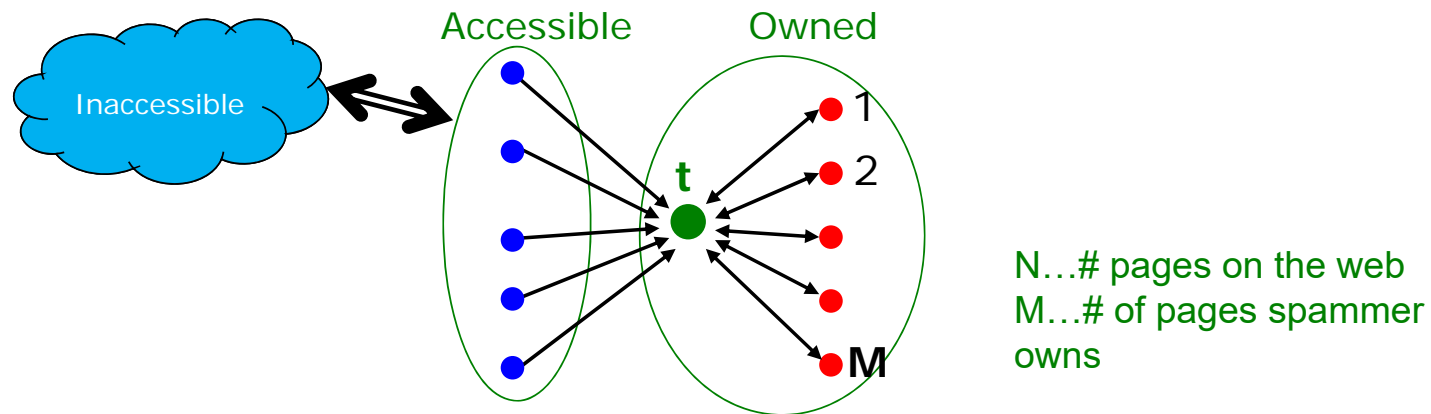
$$= x + \beta^2 y + \frac{\beta(1-\beta)M}{N} + \frac{1-\beta}{N}$$

Very small; ignore  
 Now we solve for  $y$

- $y = \frac{x}{1-\beta^2} + c \frac{M}{N}$  where  $c = \frac{\beta}{1+\beta}$



# Analysis



- $y = \frac{x}{1-\beta^2} + c \frac{M}{N}$  where  $c = \frac{\beta}{1+\beta}$
- For  $\beta = 0.85$ ,  $1/(1-\beta^2) = 3.6$ ,  $c = 0.46$
- Multiplier effect for acquired PageRank
- By making  $M$  large, we can make  $y$  as large as we want

# TrustRank: Combating the Web Spam

# Combating Spam

## ■ Combating term spam

- Analyze text using statistical methods
- Similar to email spam filtering
- Also useful: Detecting approximate duplicate pages

## ■ Combating link spam

- **Detection and blacklisting of structures that look like spam farms**
  - Leads to another war – hiding and detecting spam farms
- **TrustRank** = topic-specific PageRank with a teleport set of **trusted pages**
  - **Example:** .edu domains, similar domains for non-US schools

# TrustRank: Idea

- **Basic principle: Approximate isolation**
  - It is rare for a “good” page to point to a “bad” (spam) page
- Sample a set of **seed pages** from the web
- Have an **oracle (human)** to identify the good pages and the spam pages in the seed set
  - **Expensive task**, so we must make seed set as small as possible

# Trust Propagation

- Call the subset of seed pages that are identified as **good** the **trusted pages**
- Perform a topic-sensitive PageRank with **teleport set = trusted pages**
  - **Propagate(传播) trust through links:**
    - Each page gets a trust value between **0** and **1**
- **Solution 1: Use a threshold value and mark all pages below the trust threshold as spam**

# Simple Model: Trust Propagation

- Set trust of each trusted page to 1
- Suppose trust of page  $p$  is  $t_p$ 
  - Page  $p$  has a set of out-links  $o_p$
- For each  $q \in o_p$ ,  $p$  confers the trust to  $q$ 
  - $\beta t_p / |o_p|$  for  $0 < \beta < 1$
- Trust is additive
  - Trust of  $p$  is the sum of the trust conferred on  $p$  by all its in-linked pages
- Note similarity to Topic-Specific PageRank
  - Within a scaling factor, TrustRank = PageRank with trusted pages as teleport set

# Why is it a good idea?

## ■ Trust attenuation:

- The degree of trust conferred by a trusted page decreases with the distance in the graph

## ■ Trust splitting:

- The larger the number of out-links from a page, the less scrutiny the page author gives each out-link
- Trust is **split** across out-links

# Picking the Seed Set

- **Two conflicting considerations:**

- Human has to inspect each seed page, so seed set must be as **small** as possible
- Must ensure every **good page** gets adequate trust rank, so need make **all good pages** reachable from seed set by short paths

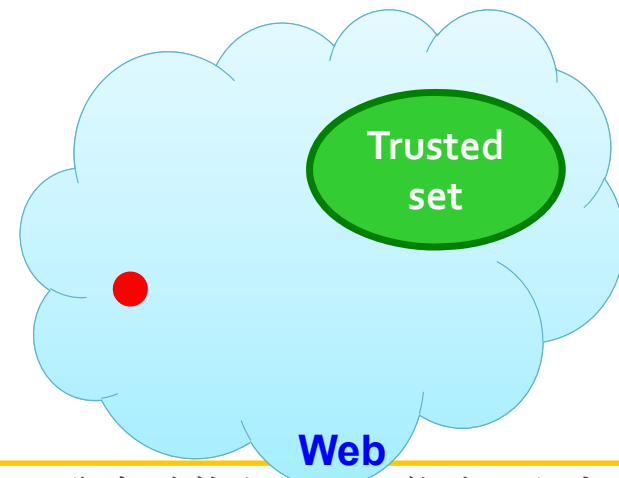


# Approaches to Picking Seed Set

- Suppose we want to pick a seed set of  $k$  pages
- **How to do that?**
- **(1) PageRank:**
  - Pick the top  $k$  pages by PageRank
  - Theory is that you can't get a bad page's rank really high
- **(2) Use trusted domains** whose membership is controlled, like .edu, .mil, .gov

# Spam Mass

- In the **TrustRank** model, we start with good pages and propagate trust
- **Complementary view:**  
What fraction of a page's PageRank comes from **spam** pages?
- In practice, we don't know all the spam pages, so we need to estimate



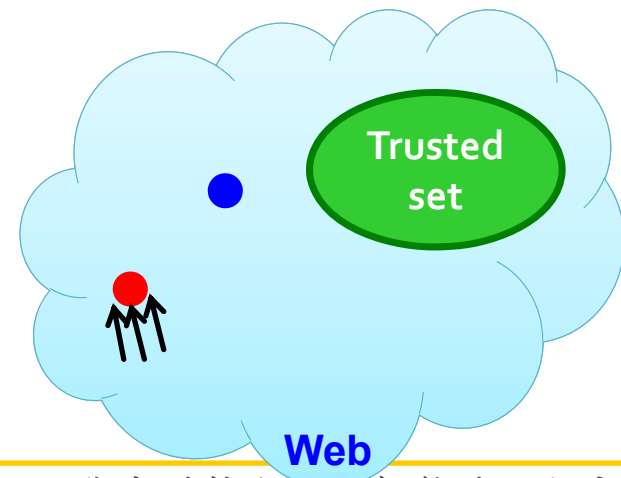
# Spam Mass Estimation

## Solution 2:

- $r_p$  = PageRank of page  $p$
- $r_p^+$  = PageRank of  $p$  with teleport into **trusted** pages only
- **Then:** What fraction of a page's PageRank comes from **spam** pages?

$$r_p^- = r_p - r_p^+$$

- **Spam mass of  $p$**   $= \frac{r_p^-}{r_p}$ 
  - Pages with high spam mass are spam.



# HITS: Hubs and Authorities

# Hubs and Authorities

- **HITS (Hypertext超文本-Induced Topic Selection)**

- Is a measure of importance of pages or documents, similar to PageRank
- Proposed at around same time as PageRank ('98)

- **Goal:** Say we want to find good newspapers

- Don't just find newspapers. Find “experts” – people who link in a coordinated way to good newspapers

- **Idea: Links as votes**

- Page is more important if it has more links
  - In-coming links? Out-going links?

# Finding newspapers

## ■ Hubs and Authorities

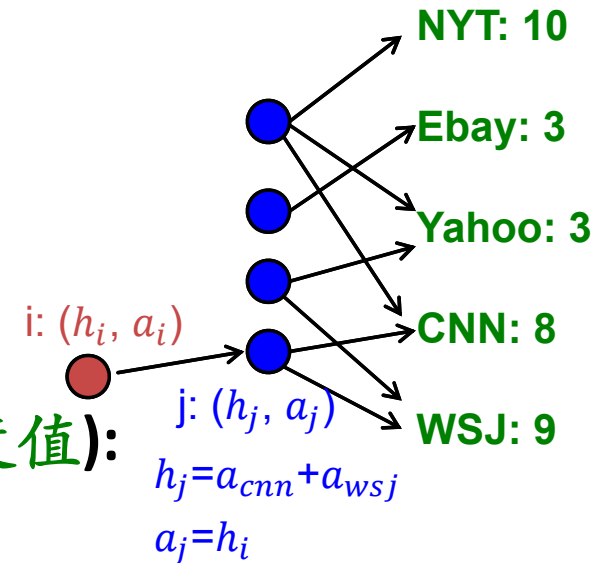
Each page has 2 scores:

### ■ Quality as an expert (**hub**, 导航度值):

- Total sum of votes of authorities pointed to

### ■ Quality as a content (**authority**, 权威度值):

- Total sum of votes coming from experts



## ■ Principle of repeated improvement

# Hubs and Authorities

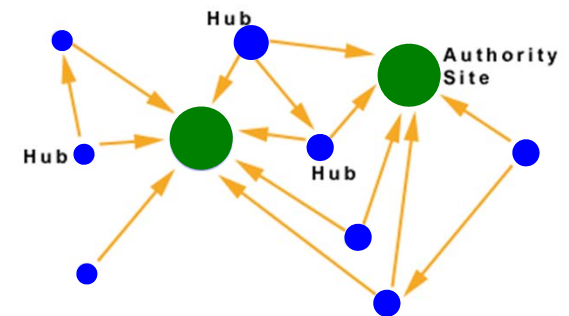
Interesting pages fall into two classes:

1. **Authorities(权威度)** are pages containing useful information

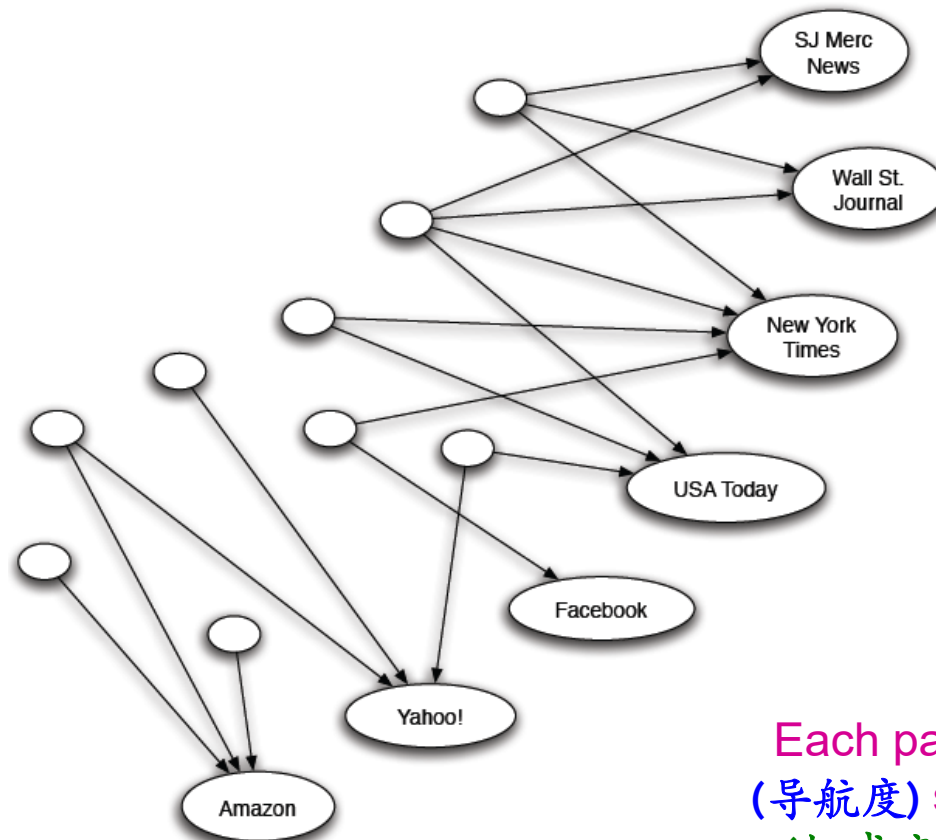
- Newspaper home pages
- Course home pages
- Home pages of auto manufacturers

2. **Hubs(导航度)** are pages that link to authorities

- List of newspapers
- Course bulletin
- List of US auto manufacturers



# Counting in-links: Authority

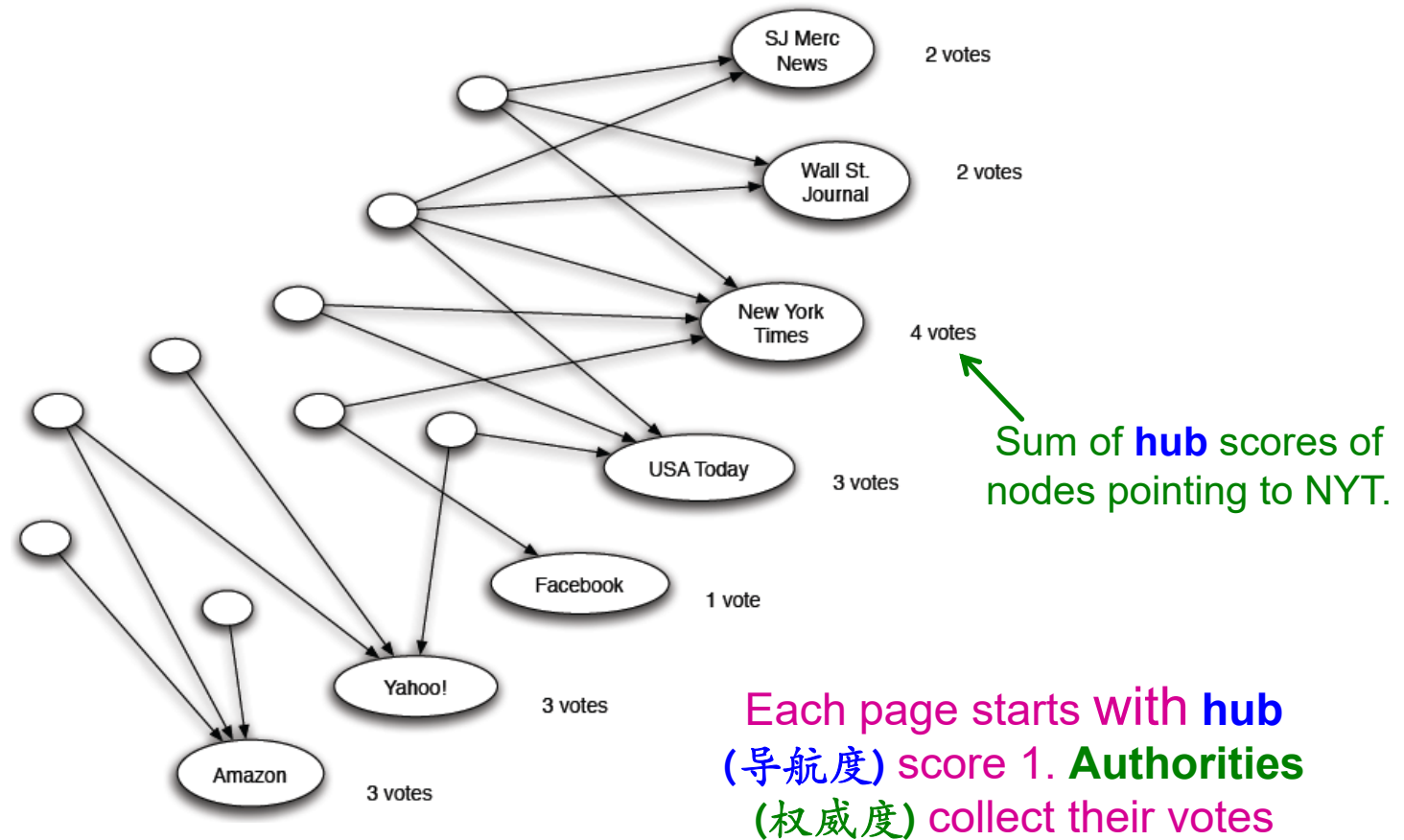


Each page starts with **hub**  
(导航度) score 1. **Authorities**  
(权威度) collect their votes

(Note this is idealized example. In reality graph is not bipartite and each page has both the hub and authority score)



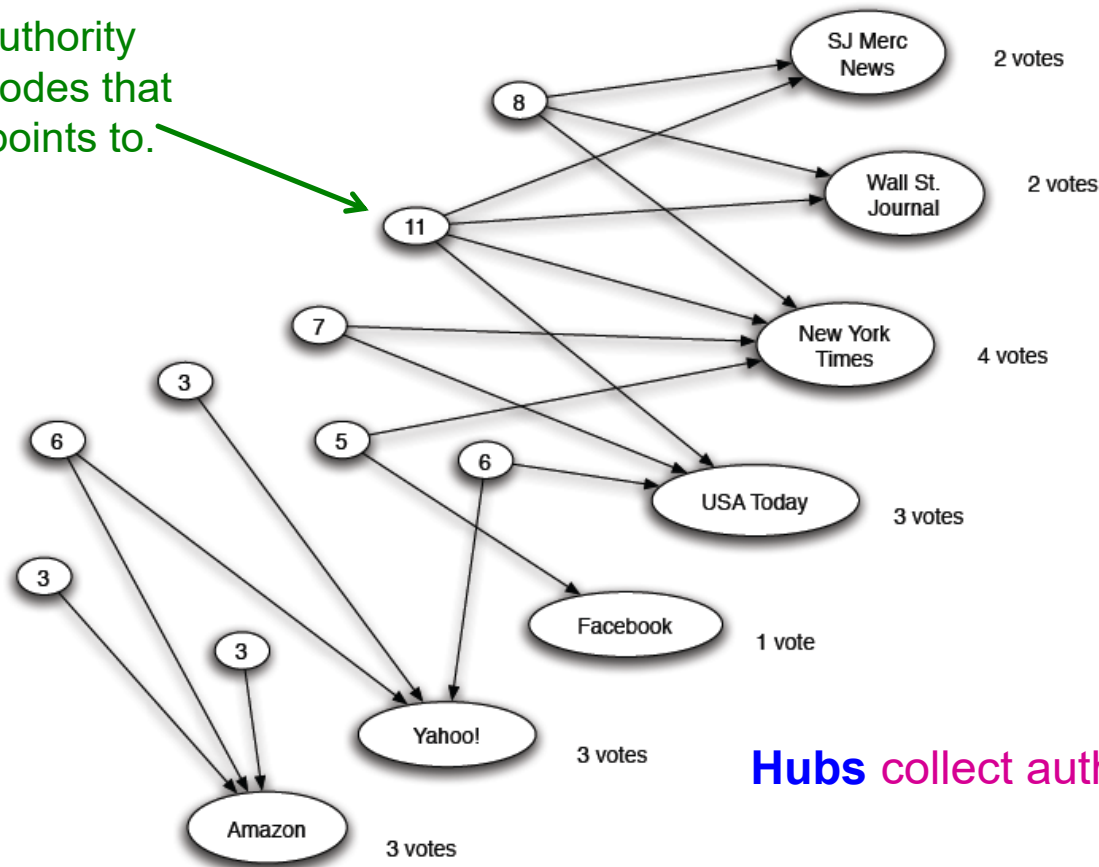
# Counting in-links: Authority



(Note this is idealized example. In reality graph is not bipartite and each page has both the hub and authority score)

# Expert Quality: Hub

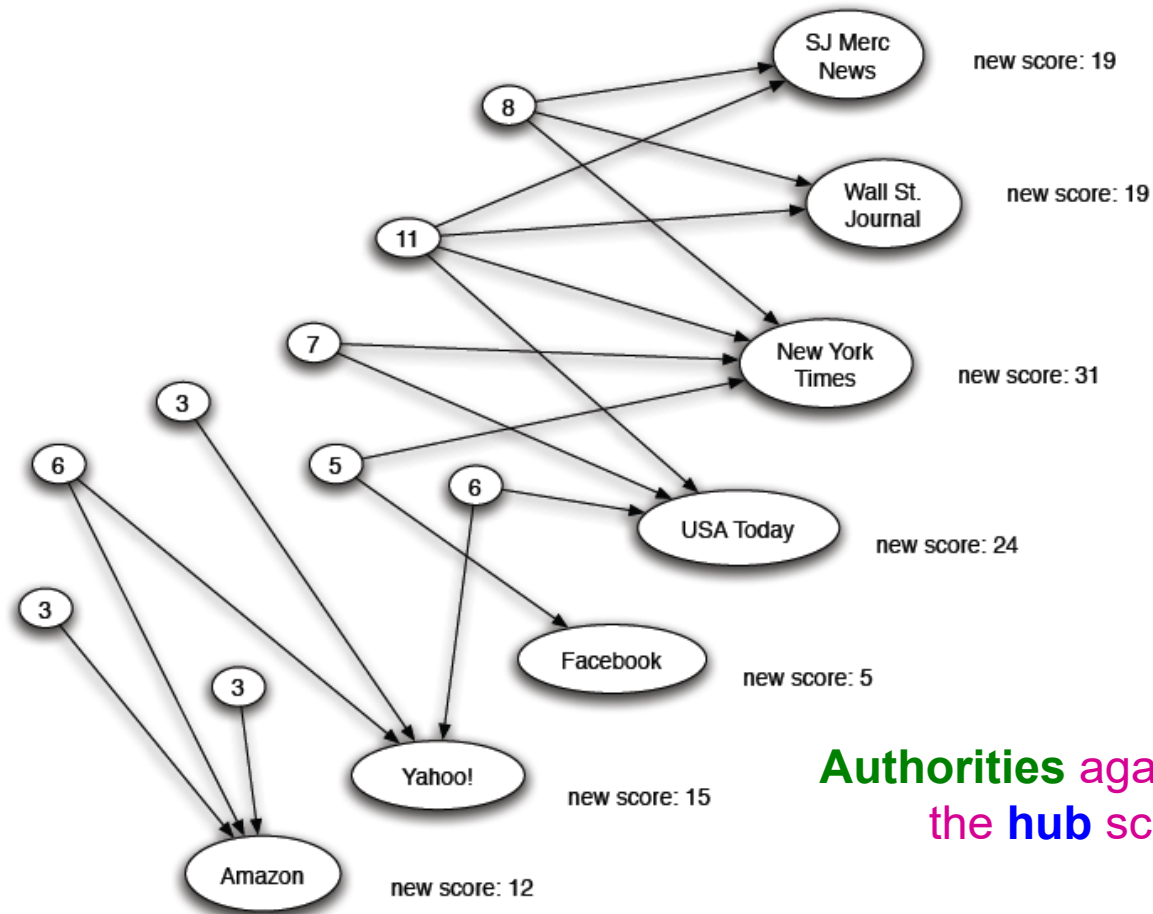
Sum of authority scores of nodes that the node points to.



**Hubs** collect authority scores

(Note this is idealized example. In reality graph is not bipartite and each page has both the hub and authority score)

# Reweighting



**Authorities** again collect  
the **hub** scores

(Note this is idealized example. In reality graph is not bipartite and each page has both the hub and authority score)