# Technologies for Web Crawling, Indexing and Search

Traian Rebedea
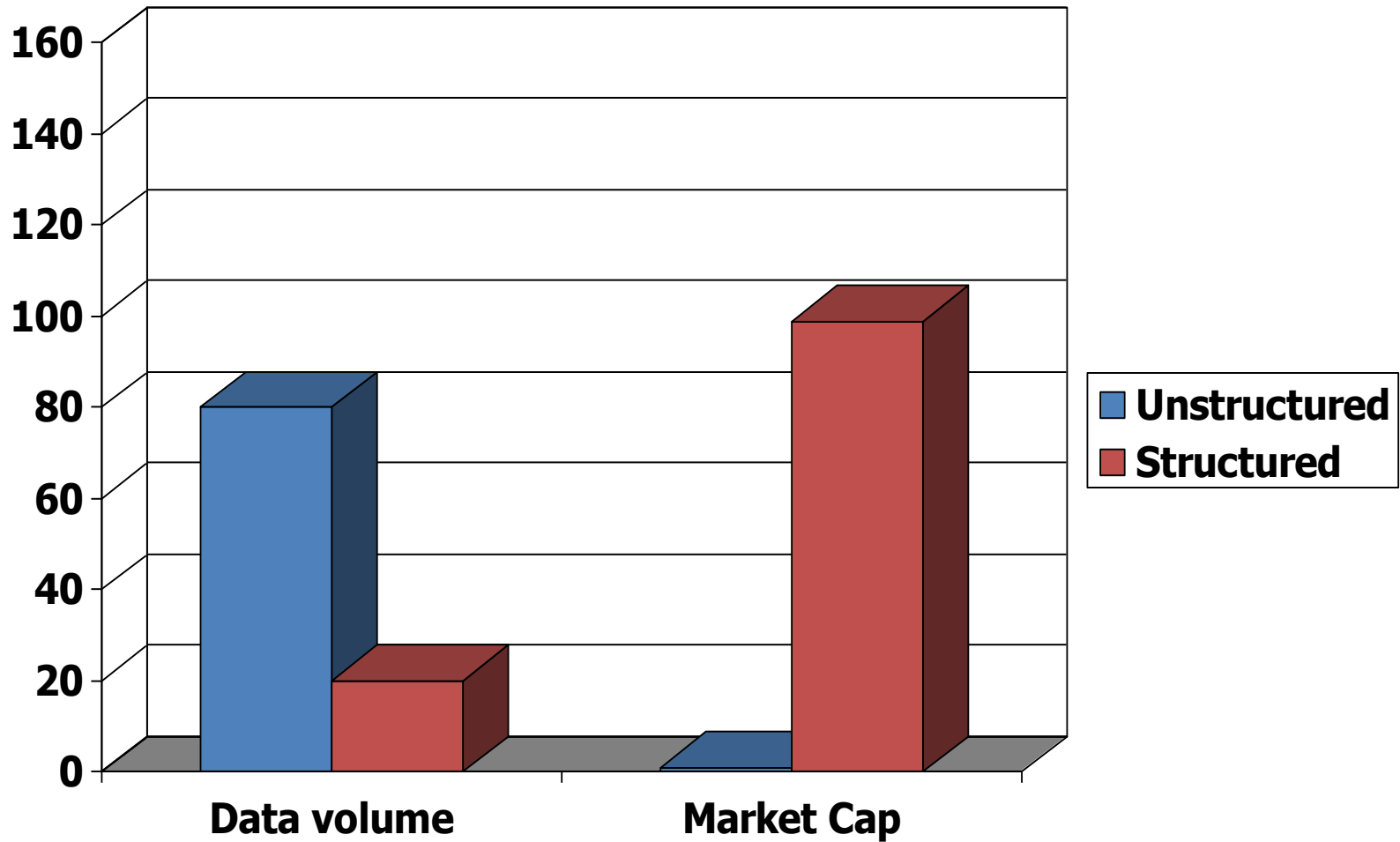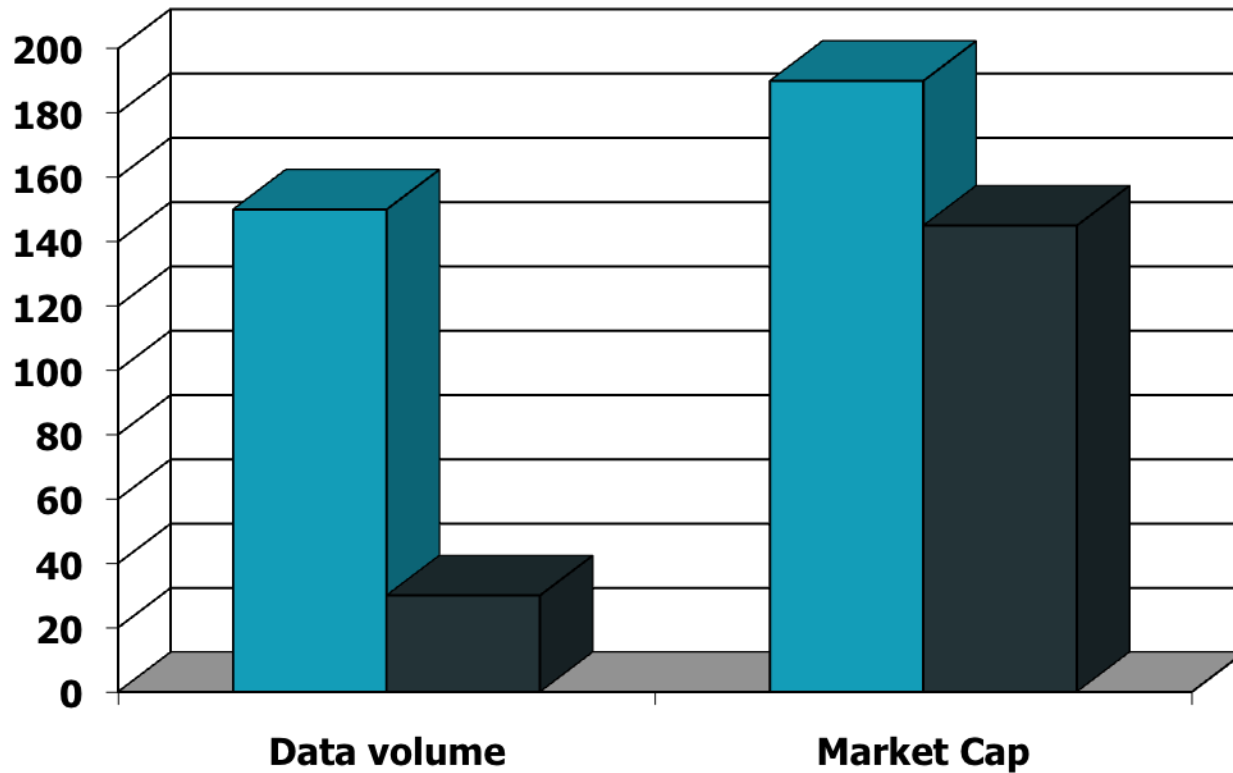
# Information Retrieval
# Search Basics

# Information Retrieval

- Information Retrieval (IR) is finding material (usually documents) of an unstructured nature (usually text) that satisfies an information need from within large collections (usually stored on computers).
  - Librarians
  - Now also in XML and DB
  - Focus on user

# Unstructured (text) vs. structured (database) data in 1996

# Unstructured (text) vs. structured (database) data in 2009

# Unstructured data in 1680

- Which plays of Shakespeare contain the words *Brutus* *AND* *Caesar* but *NOT* *Calpurnia*?

- One could grep all of Shakespeare's plays for *Brutus* and *Caesar,* then strip out lines containing *Calpurnia*?

  - Slow (for large corpora)

  - *NOT* *Calpurnia* is non-trivial

  - Other operations (e.g., find the word *Romans* near *countrymen*) not feasible

  - Ranked retrieval (best documents to return) also hard

# Solution: Term-document incidence

| | Antony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth |
|---|---|---|---|---|---|---|
| **Antony** | 1 | 1 | 0 | 0 | 0 | 1 |
| **Brutus** | 1 | 1 | 0 | 1 | 0 | 0 |
| **Caesar** | 1 | 1 | 0 | 1 | 1 | 1 |
| **Calpurnia** | 0 | 1 | 0 | 0 | 0 | 0 |
| **Cleopatra** | 1 | 0 | 0 | 0 | 0 | 0 |
| **mercy** | 1 | 0 | 1 | 1 | 1 | 1 |
| **worser** | 1 | 0 | 1 | 1 | 1 | 0 |

***Brutus** AND **Caesar** but NOT Calpurnia*

1 if play contains word, 0 otherwise

# Incidence vectors

- So we have a 0/1 vector for each term.
- To answer query: take the vectors for **Brutus, Caesar** and **Calpurnia** (complemented) ➜ bitwise *AND*.
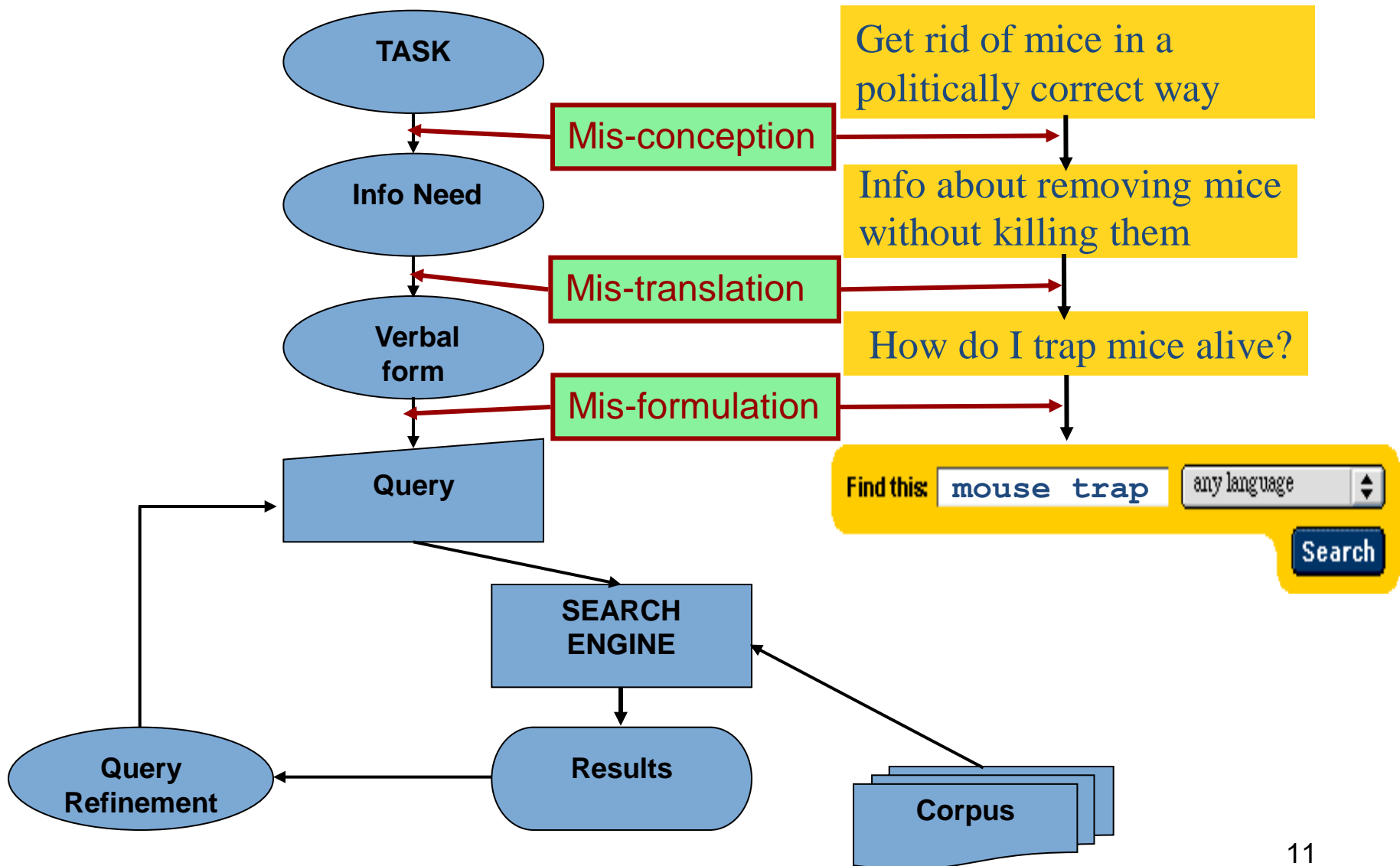- 110100 *AND* 110111 *AND* 101111 = 100100.

# Answers to query

- ## Antony and Cleopatra, Act III, Scene ii

- *Agrippa* [Aside to DOMITIUS ENOBARBUS]: Why, Enobarbus,
- When Antony found Julius **Caesar** dead,
- He cried almost to roaring; and he wept
- When at Philippi he found **Brutus** slain.


- ## Hamlet, Act III, Scene ii

- *Lord Polonius:* I did enact Julius **Caesar** I was killed i' the
- Capitol; **Brutus** killed me.

# Basic assumptions of Information Retrieval

- Corpus: Fixed document collection

- Goal: Retrieve documents with information that is <u>relevant</u> to user's information need and helps him complete a task

# The classic search model



TASK

Get rid of mice in a politically correct way

Mis-conception

Info Need

Info about removing mice without killing them

Mis-translation

Verbal form

How do I trap mice alive?

Mis-formulation

Query

Find this: mouse trap    any language    Search

SEARCH ENGINE

Results

Query Refinement

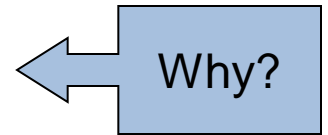Corpus

11

# How good are the retrieved docs?

- *Precision* : Fraction of retrieved docs that are relevant to user's information need

- *Recall* : Fraction of relevant docs in corpus that are retrieved

- More precise definitions and measurements to follow in later lectures

# Bigger corpora

- Consider $N$ = 1M documents, each with about 1K terms.

- Avg. 6 bytes/term incl. spaces/punctuation (EN)

  - 6GB of data in the documents.

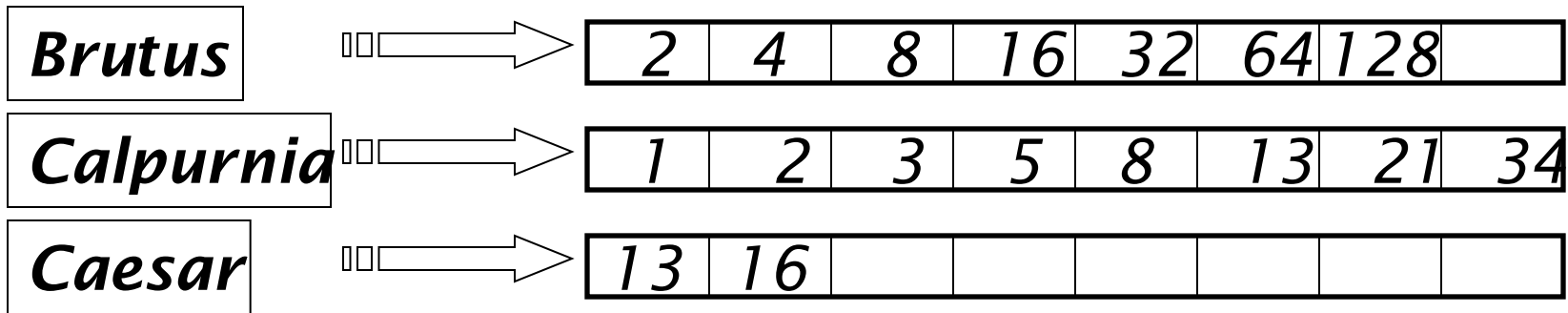- Say there are $m$ = 500K _distinct_  terms among these.

# Can't build the matrix

- 500K x 1M matrix has half-a-trillion 0's and 1's.

- But it has no more than one billion 1's.
  - matrix is extremely sparse.

  Why?

- What's a better representation?
  - We only record the 1 positions.

# Inverted index

- For each term *T*, we must store a list of all documents that contain *T*.

- Do we use an array or a list for this?

| Brutus | | 2 | 4 | 8 | 16 | 32 | 64 | 128 | |
|---|---|---|---|---|---|---|---|---|---|

| Calpurnia | | 1 | 2 | 3 | 5 | 8 | 13 | 21 | 34 |
|---|---|---|---|---|---|---|---|---|---|

| Caesar | | 13 | 16 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|

*What happens if the word **Caesar** is added to document 14?*

# Inverted index

- Linked lists generally preferred to arrays
  - Dynamic space allocation
  - Insertion of terms into documents easy
  - Space overhead of pointers

Posting

| Brutus | | 2 | 4 | 8 | 16 | 32 | 64 | 128 |

| Calpurnia | | 1 | 2 | 3 | 5 | 8 | 13 | 21 | 34 |

| Caesar | | 13 | 16 |

Dictionary

Postings lists

*Sorted by docID (more later on why).*

# Inverted index construction

**Documents to be indexed.**



Friends, Romans, countrymen.

**Token stream.**

| Friends | Romans | Countrymen |

**More on these later.**

Tokenizer

Linguistic modules

**Modified tokens.**

| friend | roman | countryman |

Indexer

**Inverted index.**

friend → 2 → 4 →

roman → 1 → 2 →

countryman → 13 → 16

# Indexer steps

- Sequence of (Modified token, Document ID) pairs.

## Doc 1

I did enact Julius
Caesar I was killed
i' the Capitol;
Brutus killed me.

## Doc 2

So let it be with
Caesar. The noble
Brutus hath told you
Caesar was ambitious

| Term | Doc # |
|------|-------|
| I | 1 |
| did | 1 |
| enact | 1 |
| julius | 1 |
| caesar | 1 |
| I | 1 |
| was | 1 |
| killed | 1 |
| i' | 1 |
| the | 1 |
| capitol | 1 |
| brutus | 1 |
| killed | 1 |
| me | 1 |
| so | 2 |
| let | 2 |
| it | 2 |
| be | 2 |
| with | 2 |
| caesar | 2 |
| the | 2 |
| noble | 2 |
| brutus | 2 |
| hath | 2 |
| told | 2 |
| you | 2 |
| caesar | 2 |
| was | 2 |
| ambitious | 2 |
|  |  |
|  |  |
|  |  |

- Sort by terms.

**Core indexing step.**

| Term | Doc # |
|------|-------|
| I | 1 |
| did | 1 |
| enact | 1 |
| julius | 1 |
| caesar | 1 |
| I | 1 |
| was | 1 |
| killed | 1 |
| i' | 1 |
| the | 1 |
| capitol | 1 |
| brutus | 1 |
| killed | 1 |
| me | 1 |
| so | 2 |
| let | 2 |
| it | 2 |
| be | 2 |
| with | 2 |
| caesar | 2 |
| the | 2 |
| noble | 2 |
| brutus | 2 |
| hath | 2 |
| told | 2 |
| you | 2 |
| caesar | 2 |
| was | 2 |
| ambitious | 2 |
| | |
| | |
| | |

➡

| Term | Doc # |
|------|-------|
| ambitious | 2 |
| be | 2 |
| brutus | 1 |
| brutus | 2 |
| capitol | 1 |
| caesar | 1 |
| caesar | 2 |
| caesar | 2 |
| did | 1 |
| enact | 1 |
| hath | 1 |
| I | 1 |
| I | 1 |
| i' | 1 |
| it | 2 |
| julius | 1 |
| killed | 1 |
| killed | 1 |
| let | 2 |
| me | 1 |
| noble | 2 |
| so | 2 |
| the | 1 |
| the | 2 |
| told | 2 |
| you | 2 |
| was | 1 |
| was | 2 |
| with | 2 |
| | |
| | |

# Indexer steps: Dictionary & Postings

- Multiple term entries in a single document are merged.
- Split into Dictionary and Postings
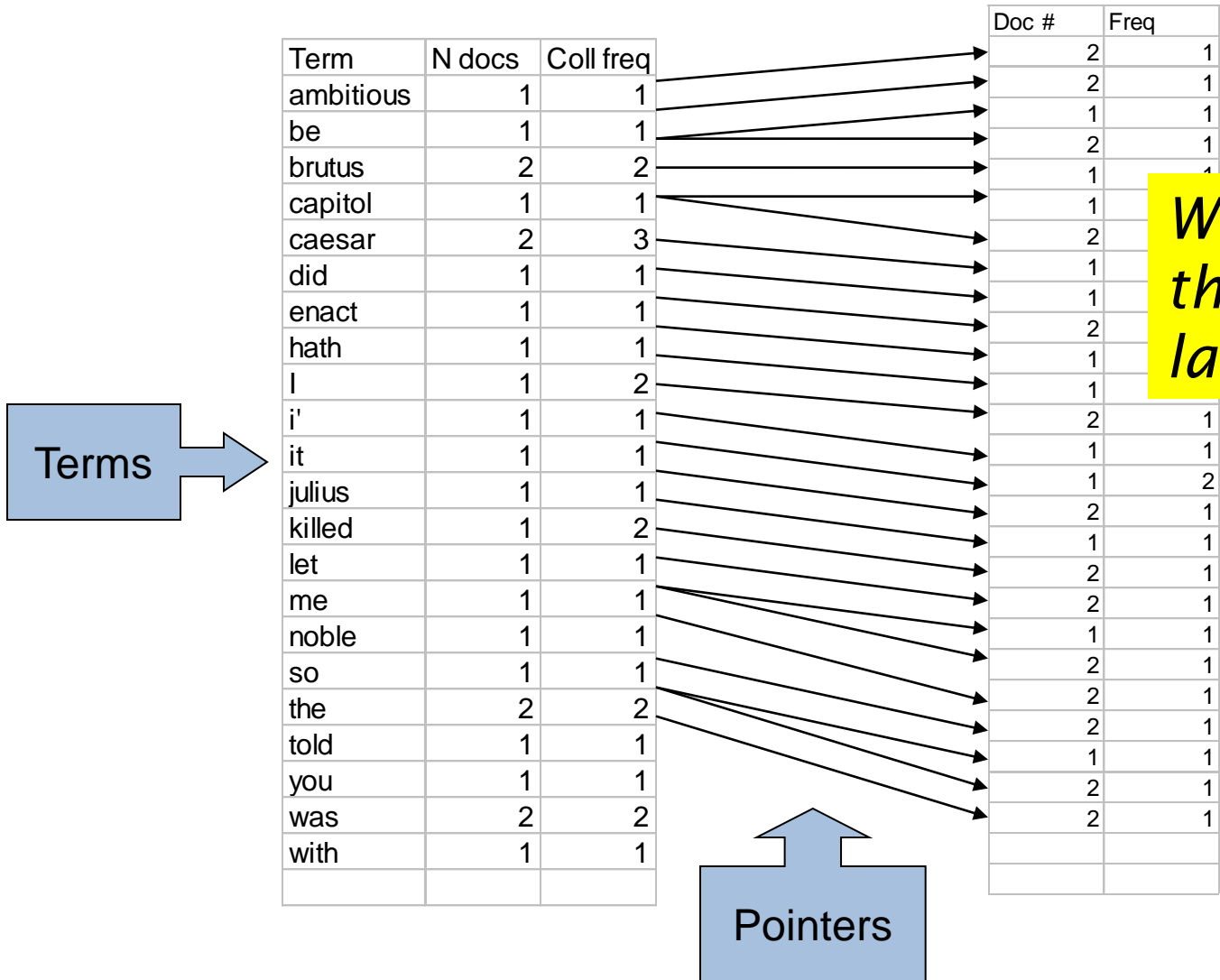- Doc. frequency information is added.

Why frequency?
Will discuss later.

| Term | docID |
|---|---|
| ambitious | 2 |
| be | 2 |
| brutus | 1 |
| brutus | 2 |
| capitol | 1 |
| caesar | 1 |
| caesar | 2 |
| caesar | 2 |
| did | 1 |
| enact | 1 |
| hath | 1 |
| I | 1 |
| I | 1 |
| i' | 1 |
| it | 2 |
| julius | 1 |
| killed | 1 |
| killed | 1 |
| let | 2 |
| me | 1 |
| noble | 2 |
| so | 2 |
| the | 1 |
| the | 2 |
| told | 2 |
| you | 2 |
| was | 1 |
| was | 2 |
| with | 2 |
| | |
| | |
| | |

| term | doc. freq. | → | postings lists |
|---|---|---|---|
| ambitious | 1 | → | 2 |
| be | 1 | → | 2 |
| brutus | 2 | → | 1 → 2 |
| capitol | 1 | → | 1 |
| caesar | 2 | → | 1 → 2 |
| did | 1 | → | 1 |
| enact | 1 | → | 1 |
| hath | 1 | → | 2 |
| i | 1 | → | 1 |
| i' | 1 | → | 1 |
| it | 1 | → | 2 |
| julius | 1 | → | 1 |
| killed | 1 | → | 1 |
| let | 1 | → | 2 |
| me | 1 | → | 1 |
| noble | 1 | → | 2 |
| so | 1 | → | 2 |
| the | 2 | → | 1 → 2 |
| told | 1 | → | 2 |
| you | 1 | → | 2 |
| was | 2 | → | 1 → 2 |
| with | 1 | → | 2 |

# Where do we pay in storage?

| Term | N docs | Coll freq |
|------|--------|-----------|
| ambitious | 1 | 1 |
| be | 1 | 1 |
| brutus | 2 | 2 |
| capitol | 1 | 1 |
| caesar | 2 | 3 |
| did | 1 | 1 |
| enact | 1 | 1 |
| hath | 1 | 1 |
| I | 1 | 2 |
| i' | 1 | 1 |
| it | 1 | 1 |
| julius | 1 | 1 |
| killed | 1 | 2 |
| let | 1 | 1 |
| me | 1 | 1 |
| noble | 1 | 1 |
| so | 1 | 1 |
| the | 2 | 2 |
| told | 1 | 1 |
| you | 1 | 1 |
| was | 2 | 2 |
| with | 1 | 1 |
| | | |

| Doc # | Freq |
|-------|------|
| 2 | 1 |
| 2 | 1 |
| 1 | 1 |
| 2 | 1 |
| 1 | 1 |
| 1 | |
| 2 | |
| 1 | |
| 1 | |
| 2 | |
| 1 | |
| 1 | |
| 2 | 1 |
| 1 | 1 |
| 1 | 2 |
| 2 | 1 |
| 1 | 1 |
| 2 | 1 |
| 2 | 1 |
| 1 | 1 |
| 2 | 1 |
| 2 | 1 |
| 2 | 1 |
| 1 | 1 |
| 2 | 1 |
| 2 | 1 |
| | |
| | |

Terms

Pointers

*Will quantify the storage, later.*

21

# The index we just built

- How do we process a query?

# Query processing: AND

- Consider processing the query:

  ***Brutus*** *AND* ***Caesar***

  - Locate ***Brutus*** in the Dictionary;

    - Retrieve its postings.
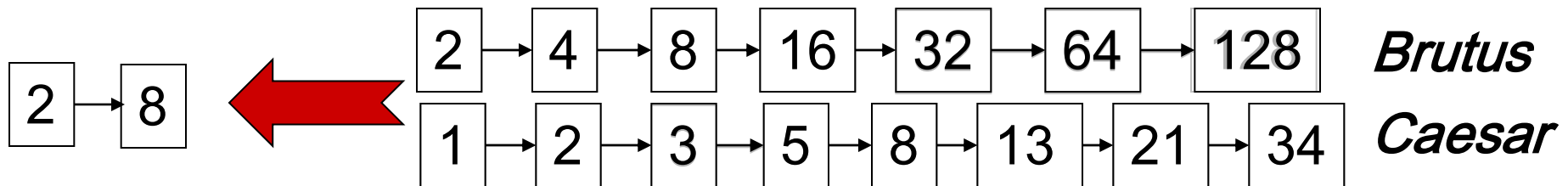
  - Locate *Caesar* in the Dictionary;

    - Retrieve its postings.

  - "Merge" the two postings:

| 2 | 4 | 8 | 16 | 32 | 64 | 128 | *Brutus* |
|---|---|---|----|----|----|-----|----------|
| 1 | 2 | 3 | 5 | 8 | 13 | 21 | 34 | *Caesar* |

# The merge

- Walk through the two postings simultaneously, in time linear in the total number of postings entries

| 2 | 8 | ← | 2 | 4 | 8 | 16 | 32 | 64 | 128 | *Brutus* |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | 1 | 2 | 3 | 5 | 8 | 13 | 21 | 34 | *Caesar* |

*If the list lengths are x and y, the merge takes O(x+y) operations.*
*<u>Crucial</u>: postings sorted by docID.*

# Intersecting two postings lists (a "merge" algorithm)

$\text{INTERSECT}(p_1, p_2)$

1   $answer \leftarrow \langle \ \rangle$

2   **while** $p_1 \neq \text{NIL}$ and $p_2 \neq \text{NIL}$

3   **do if** $docID(p_1) = docID(p_2)$

4         **then** $\text{ADD}(answer, docID(p_1))$

5                 $p_1 \leftarrow next(p_1)$

6                 $p_2 \leftarrow next(p_2)$

7         **else if** $docID(p_1) < docID(p_2)$

8                 **then** $p_1 \leftarrow next(p_1)$

9                 **else** $p_2 \leftarrow next(p_2)$

10  **return** $answer$

# Ranked Search

# Ranked retrieval

- Thus far, our queries have all been Boolean.
  - Documents either match or don't.
- Good for expert users with precise understanding of their needs and the collection.
  - Also good for applications: Applications can easily consume 1000s of results.
- Not good for the majority of users.
  - Most users incapable of writing Boolean queries (or they are, but they think it's too much work).
  - Most users don't want to wade through 1000s of results.
    - This is particularly true of web search.

# Facts

- The average query length on current search engines is 2.4 words

- Over 40% of the user queries are single words

- About 80+% of the users look only at the first page of results, 95% look at the first two pages, almost everybody looks only at the first three

# Problem with Boolean search: feast or famine

- Boolean queries often result in either too few (=0) or too many (1000s) results.

- Query 1: "*standard user dlink 650*" → 200,000 hits

- Query 2: "*standard user dlink 650 no card found*": 0 hits

- It takes a lot of skill to come up with a query that produces a manageable number of hits.

  – AND gives too few; OR gives too many

# Ranked retrieval models

- Rather than a set of documents satisfying a query expression, in ranked retrieval models, the system returns an ordering over the (top) documents in the collection with respect to a query

- Free text queries: Rather than a query language of operators and expressions, the user's query is just one or more words in a human language

- In principle, there are two separate choices here, but in practice, ranked retrieval models have normally been associated with free text queries and vice versa

# Feast or famine: not a problem in ranked retrieval

- When a system produces a ranked result set, large result sets are not an issue
  - Indeed, the size of the result set is not an issue
  - We just show the top $k$ ( ≈ 10) results
  - We don't overwhelm the user

  - Premise: the ranking algorithm works

# Scoring as the basis of ranked retrieval

- We wish to return in order the documents most likely to be useful to the searcher

- How can we rank-order the documents in the collection with respect to a query?

- Assign a score – say in [0, 1] – to each document

- This score measures how well document and query "match".

# Query-document matching scores

- We need a way of assigning a score to a query/document pair
- Let's start with a one-term query
- If the query term does not occur in the document: score should be 0
- The more frequent the query term in the document, the higher the score (should be)
- We will look at a number of alternatives for this.

# Term-document count matrices

- Consider the number of occurrences of a term in a document:
  - Each document is a count vector in $\mathbb{N}^v$: a column below

| | Antony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth |
|---|---|---|---|---|---|---|
| **Antony** | 157 | 73 | 0 | 0 | 0 | 0 |
| **Brutus** | 4 | 157 | 0 | 1 | 0 | 0 |
| **Caesar** | 232 | 227 | 0 | 2 | 1 | 1 |
| **Calpurnia** | 0 | 10 | 0 | 0 | 0 | 0 |
| **Cleopatra** | 57 | 0 | 0 | 0 | 0 | 0 |
| **mercy** | 2 | 0 | 3 | 5 | 5 | 1 |
| **worser** | 2 | 0 | 1 | 1 | 1 | 0 |

# *Bag of words* model

- Vector representation doesn't consider the ordering of words in a document
- *John is quicker than Mary* and *Mary is quicker than John* have the same vectors
- This is called the <u>bag of words</u> model.
- In a sense, this is a step back: The positional index was able to distinguish these two documents.
- We will look at "recovering" positional information later in this course.
- For now: bag of words model

# Term frequency tf

- The term frequency $\text{tf}_{t,d}$ of term $t$ in document $d$ is defined as the number of times that $t$ occurs in $d$.
- We want to use tf when computing query-document match scores. But how?
- Raw term frequency is not what we want:
  - A document with 10 occurrences of the term is more relevant than a document with 1 occurrence of the term.
  - But not 10 times more relevant.
- Relevance does not increase proportionally with term frequency.

NB: frequency = count in IR

# Log-frequency weighting

- The log frequency weight of term t in d is

$$w_{t,d} = \begin{cases} 1 + \log_{10} \mathrm{tf}_{t,d}, & \text{if } \mathrm{tf}_{t,d} > 0 \\ 0, & \text{otherwise} \end{cases}$$

- $0 \to 0$, $1 \to 1$, $2 \to 1.3$, $10 \to 2$, $1000 \to 4$, etc.
- Score for a document-query pair: sum over terms $t$ in both $q$ and $d$:
- score $= \sum_{t \in q \cap d} (1 + \log \mathrm{tf}_{t,d})$

- The score is 0 if none of the query terms is present in the document.

# Document frequency

- Rare terms are more informative than frequent terms
  - Recall stop words
- Consider a term in the query that is rare in the collection (e.g., *arachnocentric*)
- A document containing this term is very likely to be relevant to the query *arachnocentric information*
- → We want a high weight for rare terms like *arachnocentric*.

# Document frequency, continued

- Frequent terms are less informative than rare terms
- Consider a query term that is frequent in the collection (e.g., *high, increase, line*)
- A document containing such a term is more likely to be relevant than a document that doesn't
- But it's not a sure indicator of relevance.
- → For frequent terms, we want high positive weights for words like *high, increase, and line*
- But lower weights than for rare terms.
- We will use document frequency (df) to capture this.

# idf weight

- df$_t$ is the <u>document</u> frequency of $t$: the number of documents that contain $t$
  - df$_t$ is an inverse measure of the informativeness of $t$
  - df$_t \leq N$
- We define the idf (inverse document frequency) of $t$ by

  - We use log ($N$/df$_t$) instead of $N$/df$_t$ to "dampen" the effect of idf.

$$\mathrm{idf}_t = \log_{10}(N/\mathrm{df}_t)$$

Will turn out the base of the log is immaterial.

# idf example, suppose $N$ = 1 million

| term | $df_t$ | $idf_t$ |
|------|--------|---------|
| calpurnia | 1 | |
| animal | 100 | |
| sunday | 1,000 | |
| fly | 10,000 | |
| under | 100,000 | |
| the | 1,000,000 | |

$$\text{idf}_t = \log_{10}(N/\text{df}_t)$$

There is one idf value for each term $t$ in a collection.

# Effect of idf on ranking

- Does idf have an effect on ranking for one-term queries, like
  - iPhone
- idf has no effect on ranking one term queries
  - idf affects the ranking of documents for queries with at least two terms
  - For the query capricious person, idf weighting makes occurrences of capricious count for much more in the final document ranking than occurrences of person.

# tf-idf weighting

- The tf-idf weight of a term is the product of its tf weight and its idf weight.

$$\mathrm{w}_{t,d} = (1 + \log \mathrm{tf}_{t,d}) \times \log_{10}(N / \mathrm{df}_t)$$

- Best known weighting scheme in information retrieval
  - Note: the "-" in tf-idf is a hyphen, not a minus sign!
  - Alternative names: tf.idf, tf x idf
- Increases with the number of occurrences within a document
- Increases with the rarity of the term in the collection

# Final ranking of documents for a query

$$\mathrm{Score}(q,d) = \sum_{t \in q \cap d} \mathrm{tf.idf}_{t,d}$$

# Exercise 1

- Which is the ranking for the following example? (python code)
- Query: "haina cine departe"
- Document collection:
  - *D1 = "Cine împarte, parte își face"*
  - *D2 = "Cine se scoală de dimineață, departe ajunge"*
  - *D3 = "Așchia nu sare departe de trunchi"*
  - *D4 = "Omul face haina și nu haina pe om"*
  - *D5 = "Cămașa e mai aproape de piele decât haina"*

# Web Crawling

# Basic crawler operation

- Begin with known "seed" pages
- Fetch and parse them
  - Extract URLs they point to
  - Place the extracted URLs on a queue
- Fetch each URL on the queue and repeat

# Crawling picture

URLs crawled
and parsed

Unseen Web

Seed
pages

URLs *frontier*

Web

# Simple picture – complications

- Web crawling isn't feasible with one machine
  - All of the above steps distributed
- Even non-malicious pages pose challenges
  - Latency/bandwidth to remote servers vary
  - Webmasters' stipulations
    - How "deep" should you crawl a site's URL hierarchy?
  - Site mirrors and duplicate pages
- Malicious pages
  - Spam pages
  - Spider traps – including dynamically generated
- Politeness – don't hit a server too often

# What any crawler *must* do

- Be <u>Polite</u>: Respect implicit and explicit politeness considerations for a website
  - Only crawl pages you're allowed to
  - Respect *robots.txt* (more on this shortly)
- Be <u>Robust</u>: Be immune to spider traps and other malicious behavior from web servers

# What any crawler *should* do

- Be capable of <u>distributed</u> operation: designed to run on multiple distributed machines

- Be <u>scalable</u>: designed to increase the crawl rate by adding more machines

- <u>Performance/efficiency</u>: permit full use of available processing and network resources

# What any crawler *should* do

- Fetch pages of "higher <u>quality</u>" first
- <u>Continuous</u> operation: Continue fetching fresh copies of a previously fetched page
- <u>Extensible</u>: Adapt to new data formats, protocols

# Updated crawling picture



URLs crawled
and parsed

Seed
Pages

Unseen Web

URL frontier

Crawling thread

# URL frontier

- Can include multiple pages from the same host
- Must avoid trying to fetch them all at the same time
- Must try to keep all crawling threads busy

# Explicit and implicit politeness

- Explicit politeness: specifications from webmasters on what portions of site can be crawled
    - robots.txt
- Implicit politeness: even with no specification, avoid hitting any site too often

# Robots.txt

- Protocol for giving spiders ("robots") limited access to a website, originally from 1994
  - www.robotstxt.org/wc/norobots.html
- Website announces its request on what can(not) be crawled
  - For a URL, create a file `URL/robots.txt`
  - This file specifies access restrictions

# Robots.txt example

- No robot should visit any URL starting with "/yoursite/temp/", except the robot called "searchengine":

```
User-agent: *
Disallow: /yoursite/temp/

User-agent: searchengine
Disallow:
```

# Processing steps in crawling

- Pick a URL from the frontier
- Fetch the document at the URL
- Parse the URL
  - Extract links from it to other docs (URLs)
- Check if URL has content already seen
  - If not, add to indexes
- For each extracted URL
  - Ensure it passes certain URL filter
  - Check if it is already in the frontier (duplicate URL elimination)

Which one?

E.g., only crawl .edu, obey robots.txt, etc.

# Basic crawl architecture

# DNS (Domain Name Server)

- A lookup service on the internet
  - Given a URL, retrieve its IP address
  - Service provided by a distributed set of servers – thus, lookup latencies can be high (even seconds)
- Common OS implementations of DNS lookup are *blocking*: only one outstanding request at a time
- Solutions
  - DNS caching
  - Batch DNS resolver – collects requests and sends them out together

# Parsing: URL normalization

- When a fetched document is parsed, some of the extracted links are *relative* URLs

- E.g., at http://en.wikipedia.org/wiki/Main_Page

we have a relative link to /wiki/Wikipedia:General_disclaimer which is the same as the absolute URL http://en.wikipedia.org/wiki/Wikipedia:General_disclaimer

- During parsing, must normalize (expand) such relative URLs

# Content seen?

- Duplication is widespread on the web
- If the page just fetched is already in the index, do not further process it
- This is verified using document fingerprints or shingles

# Filters and robots.txt

- <u>Filters</u> – regular expressions for URL's to be crawled/not
- Once a robots.txt file is fetched from a site, need not fetch it repeatedly
  - Doing so burns bandwidth, hits web server
- Cache robots.txt files

# Duplicate URL elimination

- For a non-continuous (one-shot) crawl, test to see if an extracted+filtered URL has already been passed to the frontier

- For a continuous crawl – see details of frontier implementation

# Practical Web Crawling

- Apache Nutch (http://nutch.apache.org/)
  - Java
  - Distributed / Hadoop
  - "Using Nutch for a one of scrape of a website is like aiming a Tank at a mouse."
- Scrapy (http://scrapy.org/)
  - Python
  - Not distributed
  - Used for "scraping", not for crawling

# Practical Web Crawling (2)

- XPath is used to select elements form a DOM (Document Object Model) created from XML / HTML documents

- Example from http://vichargrave.com/xml-parsing-with-dom-using-c/

```
<bookstore>
    <book category="cooking">
        <title lang="en">Everyday Italian</title>
        <author>Giada De Laurentis</author>
        <year>2005</year>
        <price>30.00</price>
    </book>
    <book category="children">
        <title lang="en">Harry Potter and the Half-Blood Prince</title>
        <author>J. K. Rowling</author>
        <year>2005</year>
        <price>29.99</price>
    </book>
</bookstore>
```
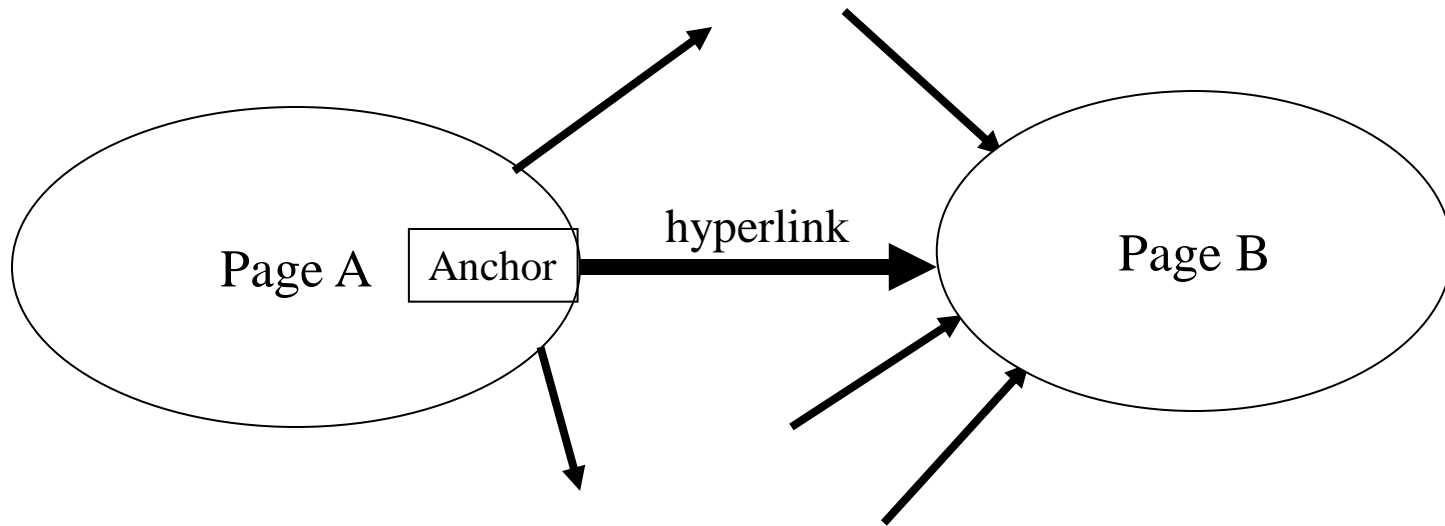
# XPath Examples

- /bookstore/book

- /bookstore/book[1]

- //book

- /bookstore/book/title[text()]

- /bookstore/book[1]/title

```xml
<bookstore>
    <book category="cooking">
        <title lang="en">Everyday Italian</title>
        <author>Giada De Laurentis</author>
        <year>2005</year>
        <price>30.00</price>
    </book>
    <book category="children">
        <title lang="en">Harry Potter and the Half-Blood Prince</title>
        <author>J. K. Rowling</author>
        <year>2005</year>
        <price>29.99</price>
    </book>
</bookstore>
```

# Exercise 2

- Crawl/scrap the news from one of the following: BBC, CNN, Reuters, NY Times, Huffington Post, Washington Post, Gandul, Hotnews, Adevarul, ...

- Install Scrapy for Python

- Read the tutorial: http://doc.scrapy.org/en/latest/intro/tutorial.html

- Write a program to extract the title and content of a news item

- Write each news item (title and content) in a different text file
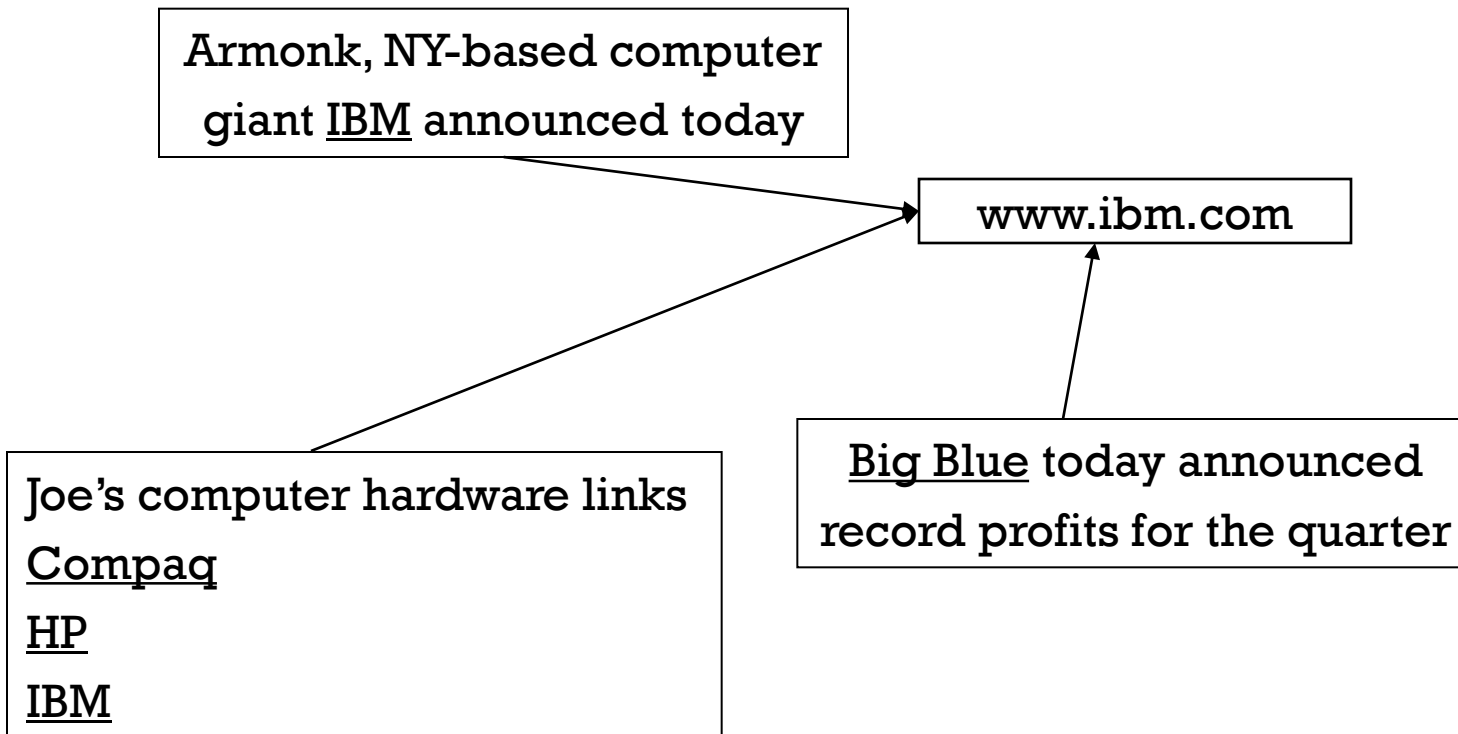
# Page Rank

# The Web as a Directed Graph



**Assumption 1:** A hyperlink between pages denotes author perceived relevance (quality signal)

**Assumption 2:** The anchor of the hyperlink describes the target page (textual context)
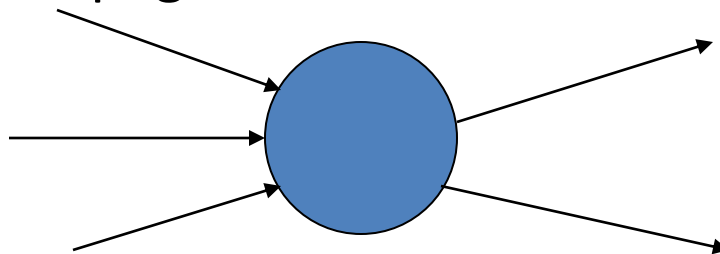
# Indexing anchor text

- When indexing a document *D*, include anchor text from links pointing to *D*.

Armonk, NY-based computer giant <u>IBM</u> announced today

www.ibm.com

Joe's computer hardware links
<u>Compaq</u>
<u>HP</u>
<u>IBM</u>

<u>Big Blue</u> today announced record profits for the quarter

# Query-independent ordering

- First generation: using link counts as simple measures of popularity.

- Two basic suggestions:
  - Undirected popularity:
    - Each page gets a score = the number of in-links plus the number of out-links (3+2=5).
  - Directed popularity:
    - Score of a page = number of its in-links (3).

# Query processing

- First retrieve all pages meeting the text query (say ***venture capital***).

- Order these by their link popularity (either variant on the previous page).

- More nuanced – use link counts as a measure of static goodness, combined with text match score
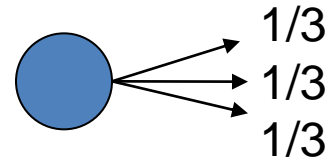
# Spamming simple popularity

- *Exercise*: How do you spam each of the following heuristics so your page gets a high score?

- Each page gets a score = the number of in-links plus the number of out-links.

- Score of a page = number of its in-links.

# Pagerank scoring

- Imagine a browser doing a random walk on web pages:
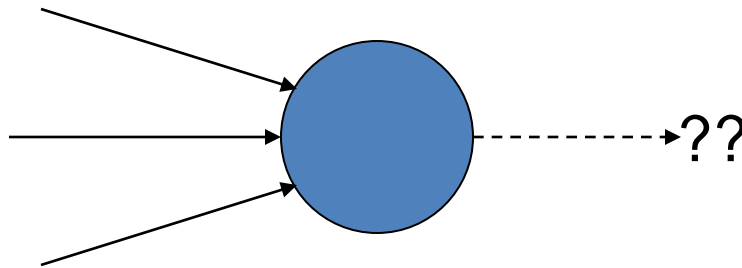
  

  1/3
  1/3
  1/3

  – Start at a random page

  – At each step, go out of the current page along one of the links on that page, equiprobably

- "In the steady state" each page has a long-term visit rate - use this as the page's score.

# Not quite enough

- The web is full of dead-ends.
    - Random walk can get stuck in dead-ends.
    - Makes no sense to talk about long-term visit rates.


??

# Teleporting

- At a dead end, jump to a random web page.
- At any non-dead end, with probability 10%, jump to a random web page.
  - With remaining probability (90%), go out on a random link.
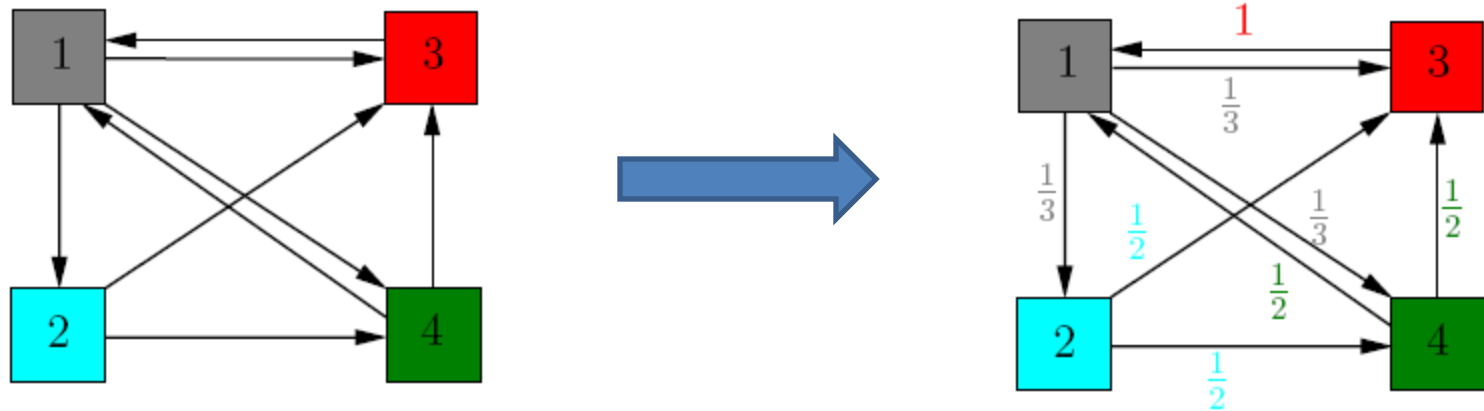  - 10% - a parameter.

# Result of teleporting

- Now cannot get stuck locally.

- There is a long-term rate at which any page is visited.

- How do we compute this visit rate?

# Web Graph

- Starting from the links, compute the weights
- This is the web graph matrix - **A**
- Example from:
  http://www.math.cornell.edu/~mec/Winter2009/RalucaRemu s/Lecture3/lecture3.html

# "Google" Matrix

- Developed by Larry Page & Sergey Brin
- Incorporates the "teleporting" solution
- Defined starting from the web graph matrix – A
- p – damping factor (usually between 0.05..0.15)

$$M = (1 - p) \cdot A + p \cdot B$$

$$B = \frac{1}{n} \cdot \begin{bmatrix} 1 & 1 & \ldots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \ldots & 1 \end{bmatrix}$$

# PageRank

- Compute the rank (importance of each page in the web graph)

- Larry Page & Sergey Brin

- Similar to citation analysis

- The rank of any page, **π**, is actually the left eigenvector of M, for the largest eigenvalue:

$$\boldsymbol{\pi\, M = \lambda\, M}$$

# Computing PageRank

- There are various methods to compute PageRank ($\pi$)

- The simplest method is called the **power (iterative) method**

- Start with an initial vector $\pi_0 = [1/n \dots 1/n]$

- Compute $\pi_{k+1} = \pi_k M$       ($k \geq 0$)

- Stop at convergence
  - Either $\pi_{k+1} = \pi_k$
  - Or $||\pi_{k+1} - \pi_k|| < \varepsilon$

# Exercise 3

- Extend the previous program in order to save the URLs and the links between these URLs

- Build the matrix A of the crawled web graph

- Build the matrix M

- Compute the PageRank of each page

- Print the URLs of the pages sorted by PageRank

# References and Further Reading

- Christopher Manning, Prabhakar Raghavan, Hinrich Schuetze: *Introduction to Information Retrieval*
- Free PDF:
  - http://nlp.stanford.edu/IR-book/information-retrieval-book.html
- Buy @ Amazon:
  - http://www.amazon.com/Introduction-Information-Retrieval-Christopher-Manning/dp/0521865719

- Most of the content in the slides has been taken from Stanford's CS276 course on Information Retrieval & Data Mining
  - http://www.stanford.edu/class/cs276/
- Many thanks to Prabhakar Raghavan for allowing the re-use of this content