

CS3245

Information Retrieval

Search Advertising,
Duplicate Detection and Revision

13

Last Time



Chapter 20

- Crawling – Obtaining documents for indexing
 - Need to be polite – `robots.txt`
 - But not everyone will return the favor – Spider Traps
 - Distributed Work (cf Distributed Indexing)

Chapter 21

- PageRank – A $G(d)$ for asymmetrically linked documents
 - Your importance hinges on who knows you



Today

Chapter 19

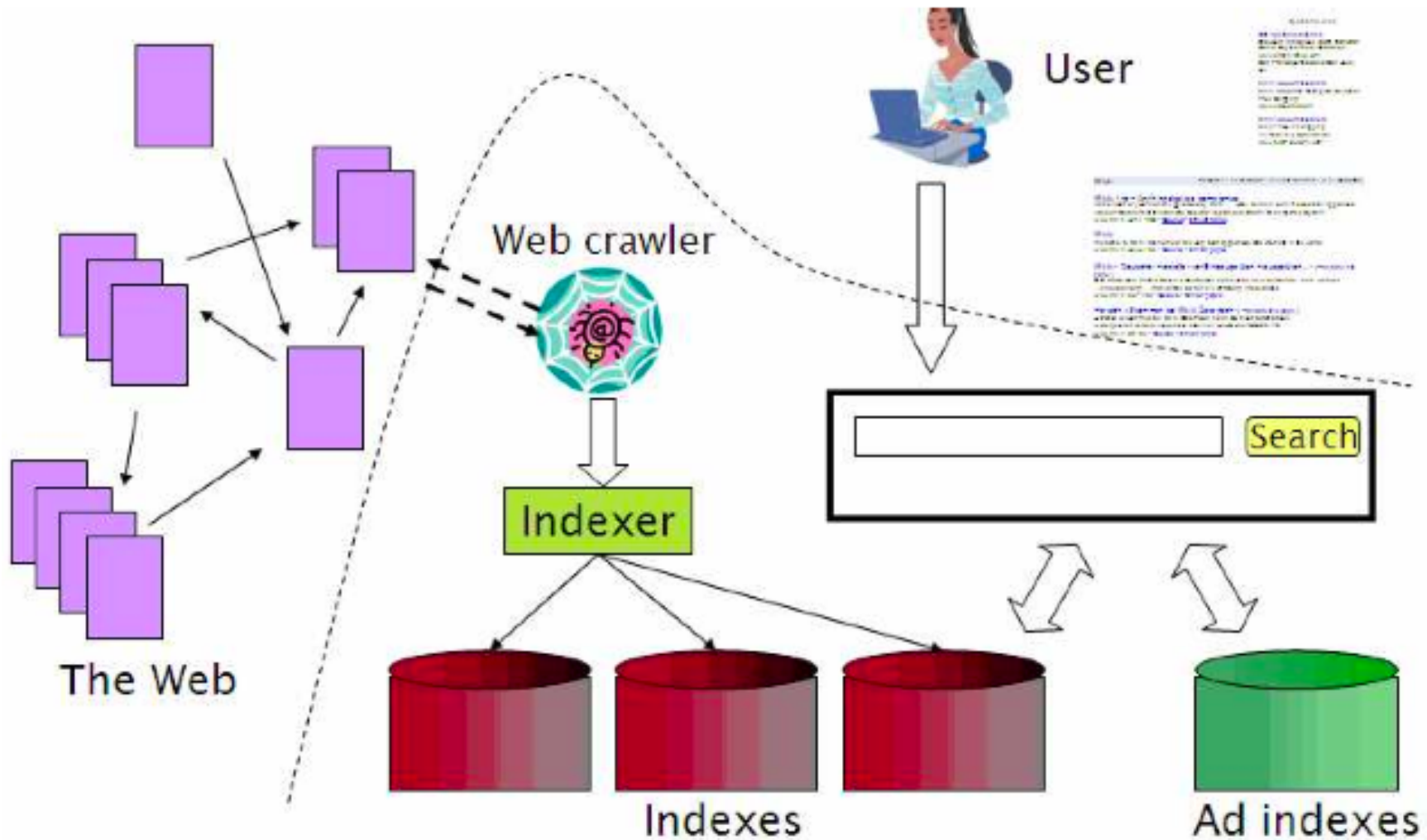
- Search Advertising
- Duplicate Detection
- Exam Format
- Revision
- Where to go from here



IR on the web vs. IR in general

- On the web, search is not just a nice feature.
 - Search is a key enabler of the web: financing, content creation, interest aggregation, etc.
 - look at search ads (Search Advertising)
- The web is a chaotic and uncoordinated document collection.
 - lots of duplicates – need to detect duplicates (Duplicate Detection)

Web search overview



Without search engines, the web wouldn't work



- Without search, content is hard to find.
- → Without search, there is no incentive to create content.
 - Why publish something if nobody will read it?
 - Why publish something if I don't get ad revenue from it?
- Somebody needs to pay for the web.
 - Servers, web infrastructure, content creation
 - A large part today is paid by search ads: Search pays for the web.

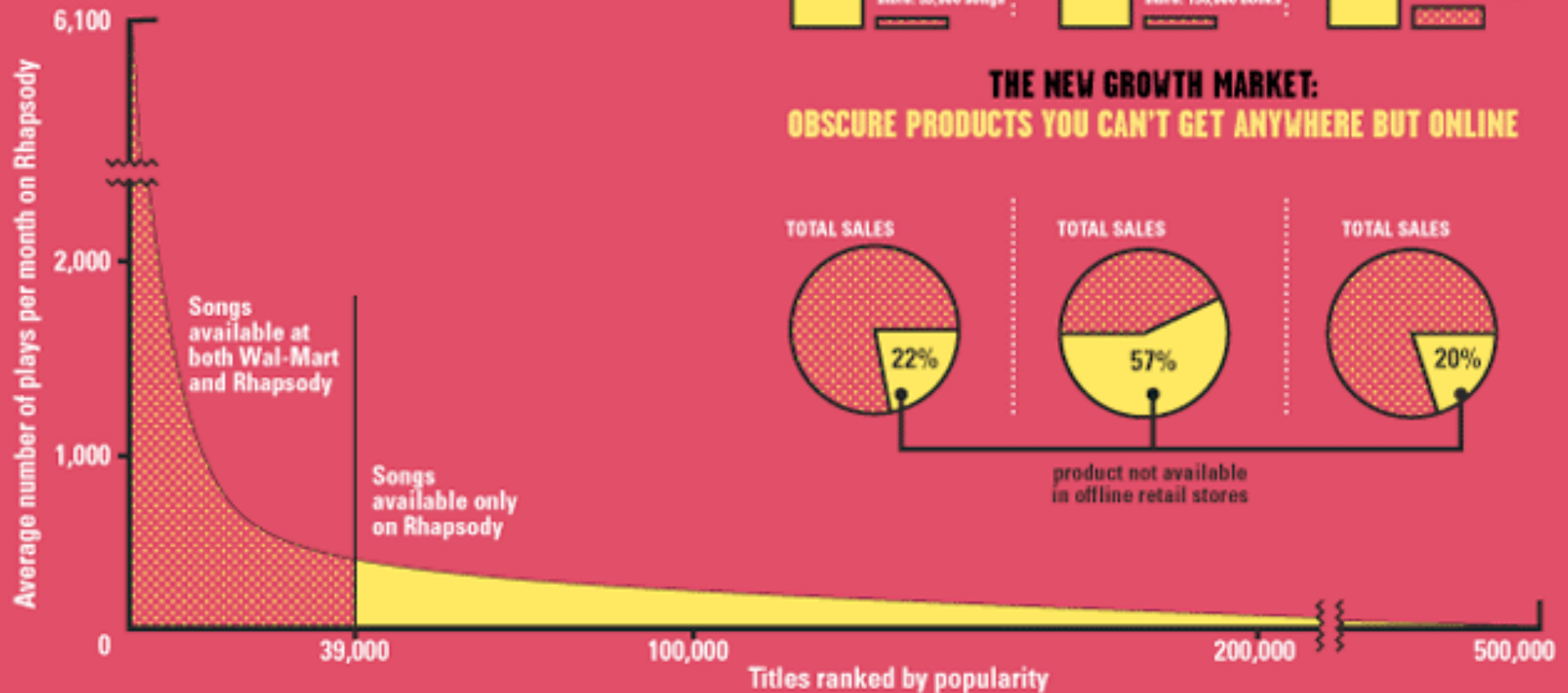
Interest aggregation



- Unique feature of the web: A small number of geographically dispersed people with similar interests can find each other.
 - Elementary school kids with hemophilia
 - People interested in translating R5R5 Scheme into relatively portable C (open source project)
 - Search engines are a key enabler for interest aggregation.

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.



Source: wired.com

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



SEARCH ADVERTISING

1st Generation of Search Ads: Goto (1996)



www.goto.com/d/search?sessionid=1AQ4214AAA6R5QF1EF3QPUQ?type=home&tm=1&Keywords=Wilmington+

Wilmington real estate.

Access 75% of all users now!
Premium Listings reach 75% of all
Internet users. [Sign up](#) for Premium
Listings today!

1. [Wilmington Real Estate - Buddy Blake](#)
Wilmington's information and real estate guide. This is your on
anything to do with Wilmington.
www.buddyblake.com (Cost to advertiser: **10.38**)
2. [Coldwell Banker Sea Coast Realty](#)
Wilmington's number one real estate company.
www.cbseacoast.com (Cost to advertiser: **10.37**)
3. [Wilmington, NC Real Estate Becky Bullard](#)
Everything you need to know about buying or selling a home c
on my Web site!
www.iwwc.net (Cost to advertiser: **10.35**)

1st Generation of Search Ads: Goto (1996)



- Buddy Blake bid the maximum (\$0.38) for this search.
- He paid \$0.38 to Goto every time somebody clicked on the link.
- Pages were simply ranked according to bid – revenue maximization for Goto.
- No separation of ads/docs. Only one result list!
- Upfront and honest. No relevance ranking, . . .
. . . but Goto did not pretend there was any.

2nd generation of search ads: Google (2000)



Web Images Maps News Shopping Gmail more [Sign in](#)

Google [Search](#) [Advanced Search](#) [Feedback](#)

Web Results 1 - 10 of about 807,000 for discount broker [\[definition\]](#). (0.12 seconds)

Discount Broker Reviews
Information on online discount brokers emphasizing rates, charges, and customer comments and complaints.
www.broker-reviews.us/ - 94k - Cached - Similar pages

Discount Broker Rankings (2003 Broker Survey) at SmartMoney.com
Discount Brokers. Rank/ Brokerage/ Minimum to Open Account, Comments, Standard Commission*, Reduced Commission, Account Fee Per Year (How to Avoid), Avg. ...
www.smartmoney.com/brokers/index.cfm?story=2004-discount-table - 121k - Cached - Similar pages

Stock Brokers | Discount Brokers | Online Brokers
Most Recommended. Top 5 Brokers headlines. 10. Don't Pay Your Broker for Free Funds May 15 at 3:39 PM. 5. Don't Discount the Discounters Apr 18 at 2:41 PM ...
www.fool.com/investing/brokers/index.aspx - 44k - Cached - Similar pages

Discount Broker
Discount Broker - Definition of Discount Broker on Investopedia - A stockbroker who carries out buy and sell orders at a reduced commission compared to a ...
www.investopedia.com/terms/d/discountbroker.asp - 31k - Cached - Similar pages

Discount Brokerage and Online Trading for Smart Stock Market ...
Online stock broker SogoTrade offers the best in discount brokerage investing. Get stock market quotes from this Internet stock trading company.
www.sogotrade.com/ - 39k - Cached - Similar pages

15 questions to ask discount brokers - MSN Money
Jan 11, 2004 ... If you're not big on hand-holding when it comes to investing, a discount broker can be an economical way to go. Just be sure to ask these ...
moneycentral.msn.com/content/Investing/StartInvesting/P68171.asp - 34k - Cached - Similar pages

Sponsored Links

Rated #1 Online Broker
No Minimums. No Inactivity Fee.
Transfer to Firsttrade for Free!
www.firsttrade.com

Discount Broker
Commission free trades for 30 days.
No maintenance fees. Sign up now.
TDAMERITRADE.com

TradeKing - Online Broker
\$4.95 per Trade, Market or Limit
SmartMoney Top Discount Broker 2001
www.TradeKing.com

Scottrade Brokerage
\$7 Trades, No Share Limit. In-Depth Research. Start Trading Online Now!
www.Scottrade.com

Stock trades \$1.95 to \$3
100 free trades, up to \$100 back for transfer costs, \$500 minimum
www.sogotrade.com

\$3.95 Online Stock Trades
Market/Limit Orders, No Share Limit and No Inactivity Fees
www.Marsco.com

INGDIRECT | ShareBuilder
Discount Broker with 2003 Best Broker Award

SogoTrade appears in search results.

SogoTrade appears in ads.

Do search engines rank advertisers higher than non-advertisers?

All major search engines claim “no”.

Do ads influence editorial content?

- Similar problem at newspapers / TV channels
- A newspaper is reluctant to publish harsh criticism of its major advertisers.
- The line often gets blurred at newspapers / on TV.
- No known case of this happening with search engines yet?
- Leads to the job of white and black hat **search engine optimization** (organic) and **search engine marketing** (paid).





How are ads ranked?

- Advertisers bid for keywords – **sale by auction**.
- Open system: Anybody can participate and bid on keywords.
- Advertisers are **only charged when somebody clicks** on your ad (i.e., CPC)

How does the auction determine an ad's **rank** and the **price paid** for the ad?

- Basis is a **second price auction**, but with twists
- For the bottom line, this is perhaps the most important research area for search engines – computational advertising.
 - Squeezing an additional fraction of **a cent** from each ad **means billions** of additional revenue for the search engine.



How are ads ranked?

- First cut: according to bid price – a la Goto
 - Bad idea: open to abuse!
 - Example: query [does my husband cheat?] → ad for divorce lawyer
 - We don't want to show nonrelevant ads.

Instead: rank based on bid price **and relevance**

- Key measure of ad relevance: clickthrough rate
 - clickthrough rate = CTR = clicks per impressions
- Result: A nonrelevant ad will be ranked low.
 - Even if this decreases search engine revenue short-term
 - Hope: Overall acceptance of the system and overall revenue is maximized if users get useful information.
- Other ranking factors: location, time of day, quality and loading speed of landing page
- The main ranking factor: the query



Google's second price auction

advertiser	bid	CTR	ad rank	rank	paid
A	\$4.00	0.01	0.04	4	(minimum)
B	\$3.00	0.03	0.09	2	\$2.68
C	\$2.00	0.06	0.12	1	\$1.51
D	\$1.00	0.08	0.08	3	\$0.51

- **bid**: maximum bid for a click by advertiser
- **CTR**: click-through rate: when an ad is displayed, what percentage of time do users click on it? **CTR is a measure of relevance.**
- **ad rank**: $\text{bid} \times \text{CTR}$: this trades off (i) how much money the advertiser is willing to pay against (ii) how relevant the ad is
- **rank**: rank in auction
- **paid**: second price auction price paid by advertiser



Google's second price auction

advertiser	bid	CTR	ad rank	rank	paid
A	\$4.00	0.01	0.04	4	(minimum)
B	\$3.00	0.03	0.09	2	\$2.68
C	\$2.00	0.06	0.12	1	\$1.51
D	\$1.00	0.08	0.08	3	\$0.51

- Second price auction: The advertiser pays the minimum amount necessary to maintain their position in the auction (plus 1 cent) – related to the Vickrey Auction
- $\text{price}_1 \times \text{CTR}_1 = \text{bid}_2 \times \text{CTR}_2$ (this will result in $\text{rank}_1 = \text{rank}_2$)
- $\text{price}_1 = \text{bid}_2 \times \text{CTR}_2 / \text{CTR}_1$
- $p_1 = \text{bid}_2 \times \text{CTR}_2 / \text{CTR}_1 = 3.00 \times 0.03 / 0.06 = 1.50$
- $p_2 = \text{bid}_3 \times \text{CTR}_3 / \text{CTR}_2 = 1.00 \times 0.08 / 0.03 = 2.67$
- $p_3 = \text{bid}_4 \times \text{CTR}_4 / \text{CTR}_3 = 4.00 \times 0.01 / 0.08 = 0.50$

Keywords with high bids



According to <http://www.cwire.org/highest-paying-search-terms/>

\$69.1	mesothelioma treatment options
\$65.9	personal injury lawyer michigan
\$62.6	student loans consolidation
\$61.4	car accident attorney los angeles
\$59.4	online car insurance quotes
\$59.4	arizona dui lawyer
\$46.4	asbestos cancer
\$40.1	home equity line of credit
\$39.8	life insurance quotes
\$39.2	refinancing
\$38.7	equity line of credit
\$38.0	lasik eye surgery new york city
\$37.0	2nd mortgage
\$35.9	free car insurance quote

Also more recently:

<http://www.wordstream.com/articles/most-expensive-keywords>

Search ads: A win-win-win?



- The **search engine** company gets revenue every time somebody clicks on an ad.
- The **user** only clicks on an ad if they are interested in the ad.
 - Search engines punish misleading and nonrelevant ads.
 - As a result, users are often satisfied with what they find after clicking on an ad.
- The **advertiser** finds new customers in a cost-effective way.

Not a win-win-win: Keyword arbitrage



- Buy a keyword on Google
 - Then redirect traffic to a third party that is paying much more than you are paying Google.
 - E.g., redirect to a page full of ads
 - This rarely makes sense for the user.
 - (Ad) spammers keep inventing new tricks.
 - The search engines need time to catch up with them.
- ➔ Adversarial Information Retrieval

Not a win-win-win: Violation of trademarks



“geico”

- During part of 2005: The search term “geico” on Google was bought by competitors.
- Geico lost this case in the United States.
- Louis Vuitton lost a similar case in Europe.

<http://support.google.com/adwordspolicy/answer/6118?rd=1>

- It's potentially misleading to users to trigger an ad off of a trademark if the user can't buy the product on the site.



DUPLICATE DETECTION

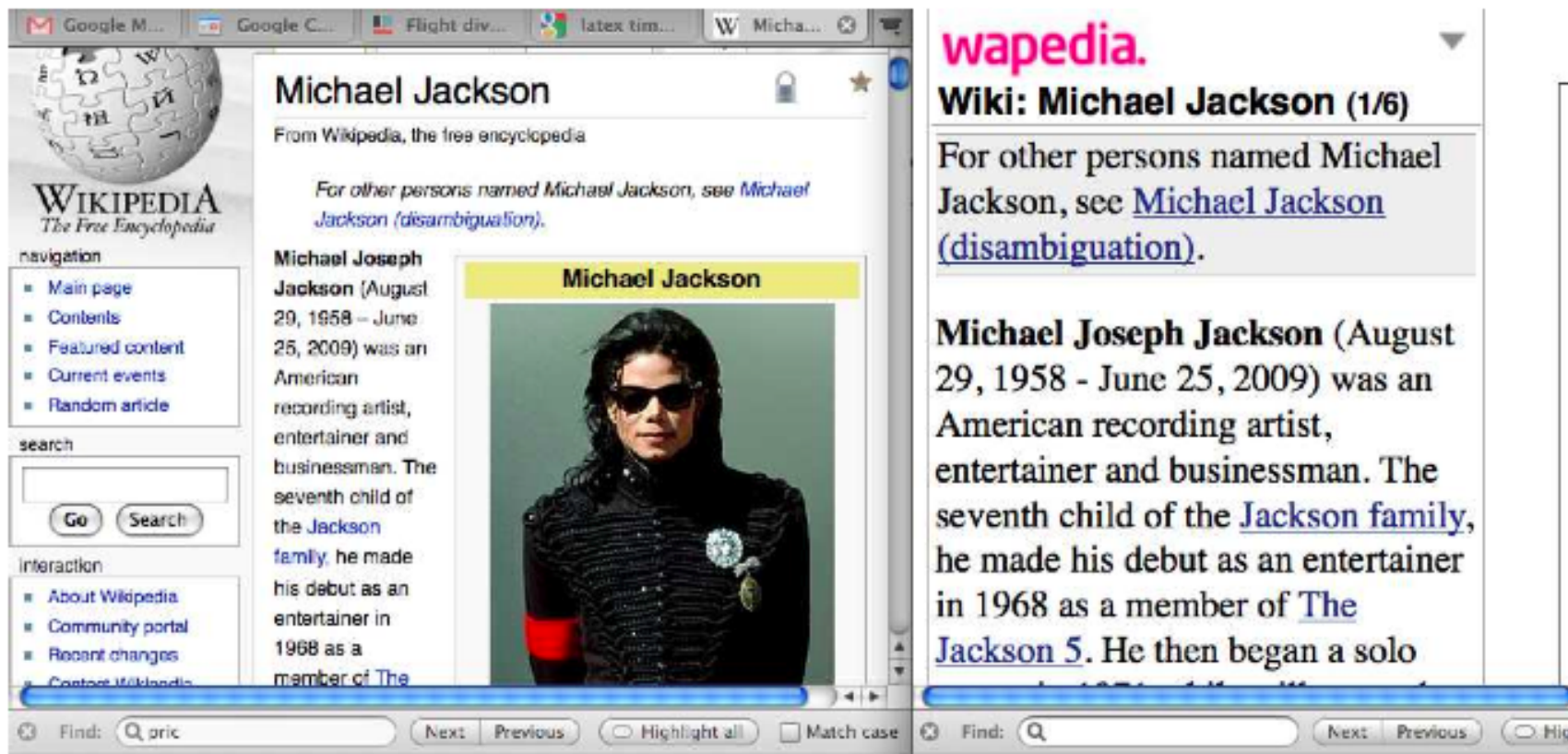


Duplicate detection

The web is full of duplicated content.

- More so than many other collections
- Exact duplicates
 - Easy to detect; use hash/fingerprint (e.g., MD5)
- Near-duplicates
 - More common on the web, difficult to eliminate
- For the user, it's annoying to get a search result with near-identical documents.
- **Marginal relevance is zero**: even a highly relevant document becomes nonrelevant if it appears below a (near-)duplicate.

Near-duplicates: Example



The image shows a side-by-side comparison of two web pages for Michael Jackson. The left page is the Wikipedia article, and the right page is a Wapedia article. Both pages contain a title, a brief biography, and a photograph of Michael Jackson. The Wapedia page has a pink header and a different layout for the navigation and search sections.

Wikipedia:

Michael Jackson
From Wikipedia, the free encyclopedia

For other persons named Michael Jackson, see [Michael Jackson \(disambiguation\)](#).

Michael Joseph Jackson (August 29, 1958 – June 25, 2009) was an American recording artist, entertainer and businessman. The seventh child of the [Jackson family](#), he made his debut as an entertainer in 1968 as a member of [The](#)

Wapedia:

wapedia.
Wiki: Michael Jackson (1/6)

For other persons named Michael Jackson, see [Michael Jackson \(disambiguation\)](#).

Michael Joseph Jackson (August 29, 1958 - June 25, 2009) was an American recording artist, entertainer and businessman. The seventh child of the [Jackson family](#), he made his debut as an entertainer in 1968 as a member of [The Jackson 5](#). He then began a solo

Detecting near-duplicates



- Compute similarity with an edit-distance measure
- We want “**syntactic**” (as opposed to **semantic**) similarity.
 - True semantic similarity (similarity in content) is too difficult to compute.
- We do not consider documents near-duplicates if they have the same content, but express it with different words.
- Use similarity threshold θ to make the call “is/isn’t a near-duplicate”.
- E.g., two documents are near-duplicates if similarity
- $> \theta = 80\%$.

Recall: Jaccard coefficient



- A commonly used measure of overlap of two sets
- Let A and B be two sets
- Jaccard coefficient:

$$\text{JACCARD}(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

$(A \neq \emptyset \text{ or } B \neq \emptyset)$

- $\text{JACCARD}(A, A) = 1$
- $\text{JACCARD}(A, B) = 0$ if $A \cap B = \emptyset$
- A and B don't have to be the same size.
- Always assigns a number between 0 and 1.



Jaccard coefficient: Example



- Three documents:
- d_1 : “Jack London traveled to Oakland”
- d_2 : “Jack London traveled to the city of Oakland”
- d_3 : “Jack traveled from Oakland to London”
- Based on shingles of size 2 (2-grams or bigrams), what are the Jaccard coefficients $J(d_1, d_2)$ and $J(d_1, d_3)$?
- $J(d_1, d_2) = 3/8 = 0.375$
- $J(d_1, d_3) = 0$
- **Note:** very sensitive to dissimilarity



A document as set of shingles

- A shingle is simply a **word n-gram**.
- Shingles are used as features to **measure syntactic similarity** of documents.
- For example, for $n = 3$, “a rose is a rose is a rose” would be represented as this set of shingles:
 - { a-rose-is, rose-is-a, is-a-rose }
- We define the similarity of two documents as the **Jaccard coefficient of their shingle sets**.



Fingerprinting

- We can map shingles to a large integer space $[1..2^m]$ (e.g., $m = 64$) by fingerprinting.
- We use s_k to refer to the shingle's fingerprint in $1..2^m$.
- This doesn't directly help us – we are just converting strings to large integers
- But it **will** help us compute an approximation to the actual Jaccard coefficient quickly



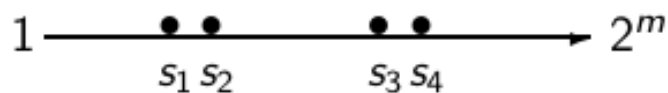
Documents as sketches

- The number of shingles per document is large, difficult to exhaustively compare
- To make it fast, we use a **sketch**, a **subset** of the shingles of a document.
- The size of a sketch is, say, $n = 200$ and is defined by a set of permutations $\pi_1 \dots \pi_{200}$.
- Each π_i is a random permutation on $1..2^m$
- The **sketch** of d is defined as:
 - $\langle \min_{s \in d} \pi_1(s), \min_{s \in d} \pi_2(s), \dots, \min_{s \in d} \pi_{200}(s) \rangle$
(a vector of 200 numbers).

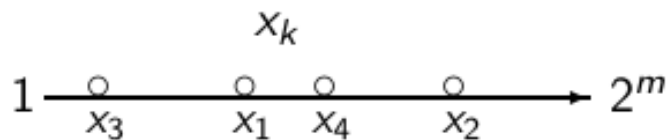
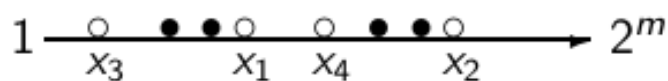
Deriving a sketch element: a permutation of the original hashes



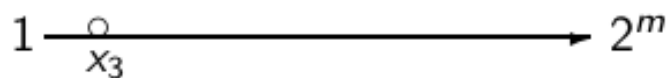
document 1: $\{s_k\}$



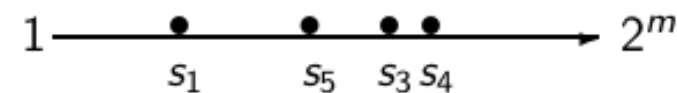
$$x_k = \pi(s_k)$$



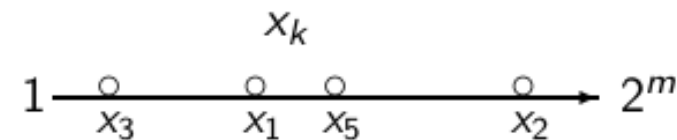
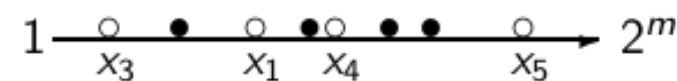
$$\min_{s_k} \pi(s_k)$$



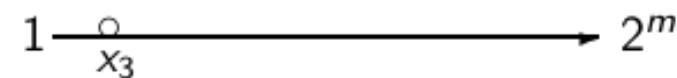
document 2: $\{s_k\}$



$$x_k = \pi(s_k)$$



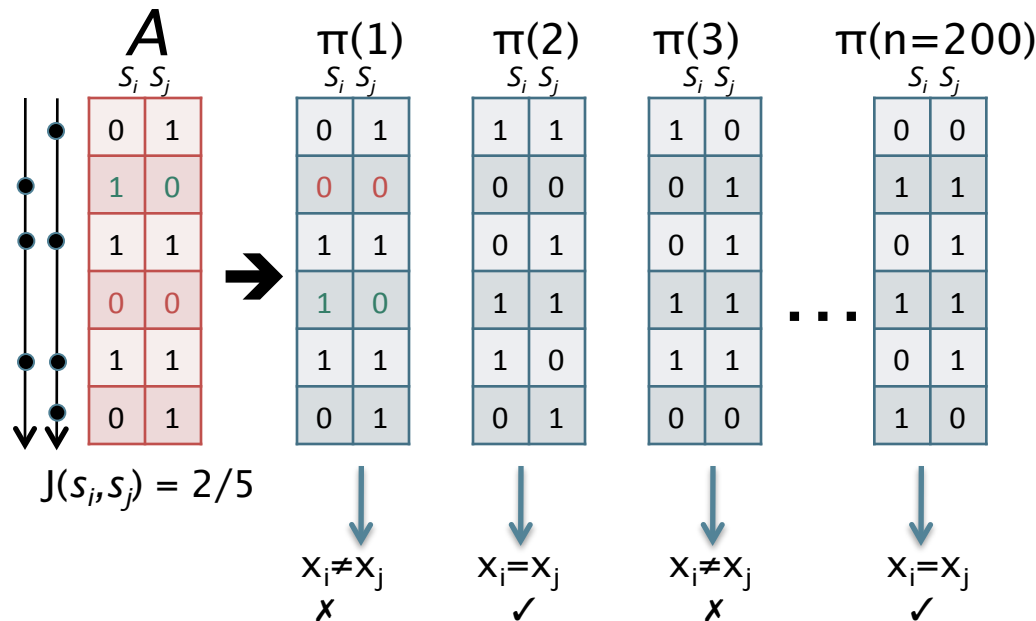
$$\min_{s_k} \pi(s_k)$$



We use $\min_{s \in d_1} \pi(s) = \min_{s \in d_2} \pi(s)$ as a test for: are d_1 and d_2 near-duplicates? In this case: permutation π says: $d_1 \approx d_2$



Proof that $J(S(d_i), s(d_j)) \cong P(x_i^\pi = x_j^\pi)$



We view a matrix A :

- 1 column per set of hashes
- Element $A_{i,j} = 1$ if element i in set S_j is present
- Permutation $\pi(n)$ is a random reordering of the rows in A
- x_k^π is the first non-zero entry in $\pi(d_k)$, i.e., first shingle present in document k

- Let $C_{00} = \#$ of rows in A where both entries are 0, define C_{11}, C_{10}, C_{01} likewise.
- $J(s_i, s_j)$ is then equivalent to $C_{11} / C_{10} + C_{01} + C_{11}$
- $P(x_i = x_j)$ then is equivalent to $C_{11} / C_{10} + C_{01} + C_{11}$



Estimating Jaccard

- Thus, the proportion of successful permutations is the Jaccard coefficient.
 - Permutation π is successful iff $\min_{s \in d_1} \pi(s) = \min_{s \in d_2} \pi(s)$
- Picking a permutation at random and outputting 1 (successful) or 0 (unsuccessful) is a Bernoulli trial.
- Estimator of probability of success: proportion of successes in n Bernoulli trials. ($n = 200$)
- Our sketch is based on a random selection of permutations.
- Thus, to compute Jaccard, count the number k of successful permutations for $\langle d_1, d_2 \rangle$ and divide by $n = 200$.
- $k/n = k/200$ estimates $J(d_1, d_2)$



Shingling: Summary

- Input: N documents
- Choose n -gram size for shingling, e.g., $n = 5$
- Pick 200 random permutations, represented as hash functions
- Compute N sketches: $200 \times N$ matrix shown on previous slide, one row per permutation, one column per document
- Compute $\frac{N \cdot (N-1)}{2}$ pairwise similarities
- Transitive closure of documents with similarity $> \theta$
- Index only one document from each equivalence class



Summary

- Search Advertising
 - A type of crowdsourcing: Ask advertisers how much they want to spend, ask searchers how relevant an ad is
 - Auction Mechanics

- Duplicate Detection
 - Represent documents as shingles
 - Calculate an approximation to the Jaccard by using random trials.



EXAM TOPICS



Exam Format

Open Book

Topics
not in order

INSTRUCTIONS TO CANDIDATES

1. This examination paper contains SIX (6) questions and comprises NINE (9) printed pages, including this page. Some questions have multiple parts.
2. It is suggested that you limit your response length to the space in the boxes provided.
3. You may use the backs of the pages as scratch paper, as they *will* be disregarded, unless specifically noted by you.
4. This is an OPEN BOOK examination. You may consult books and any other printed or handwritten materials for this test.
5. You may use pencil or other erasable medium in answering this paper.
6. The questions are presented no particular order, and specifically not by their perceived difficulty or estimated time to answer. You may want to do the questions out of order.
7. Please write your Matriculation Number below. Do not write your name.

MATRICULATION NO: _____

This portion is for examiner's use only

Question	Q1	Q2	Q3	Q4	Q5	Q6	Total
Max	20	16	23	20	10	11	100
Marks							



Exam Topics

- Anything covered in lecture (through slides)
- And corresponding sections in textbook
- Tutorial and homework essays are the exam models
 - In the past, I wrote the exam questions when writing tutorials
 - There is also a graduate level module that I've taught before (CS 5246; I taught it in 2007)
- Not responsible for sections that we didn't cover
 - If in doubt, ask on the forum

Exam Topic Distribution

- Emphasize on second half of semester, but skips some topics
- 1 or 2 questions may have calculation
 - May be time consuming but easy and straightforward
 - **Don't forget your calculator!**
- Others are thinking essay questions (cf. tutorials)

Week 1 (16 Jan) [Course overview and Language Models](#)

Week 2 (23 Jan) [Boolean Retrieval](#)

Week 3 (30 Jan) [Postings List and Choosing Terms](#)

Week 4 (6 Feb) [Dictionaries and Tolerant Retrieval](#)

Week 5 (13 Feb) [Index Compression](#)

Week 6 (20 Feb) No class in lieu of Chinese (Lunar) New Year -- Watch recorded lecture on [Index Construction](#)

Recess Week

Week 7 (6 Mar) [Scoring and the Vector Space Model](#)

Week 8 (13 Mar) NUS E-Learning Week -- Watch recorded lecture on [A Complete Search System](#)

Week 9 (20 Mar) [IR Evaluation](#)

Week 10 (27 Mar) [XML IR, Relevance Feedback and Query Expansion](#)

Week 11 (3 Apr) No class in lieu of Good Friday -- Watch recorded lecture on [Probabilistic IR and Language Model IR](#) or come to in-place make-up lecture (venue and timing TBA)

Week 12 (10 Apr) [Crawling and Link Analysis](#)

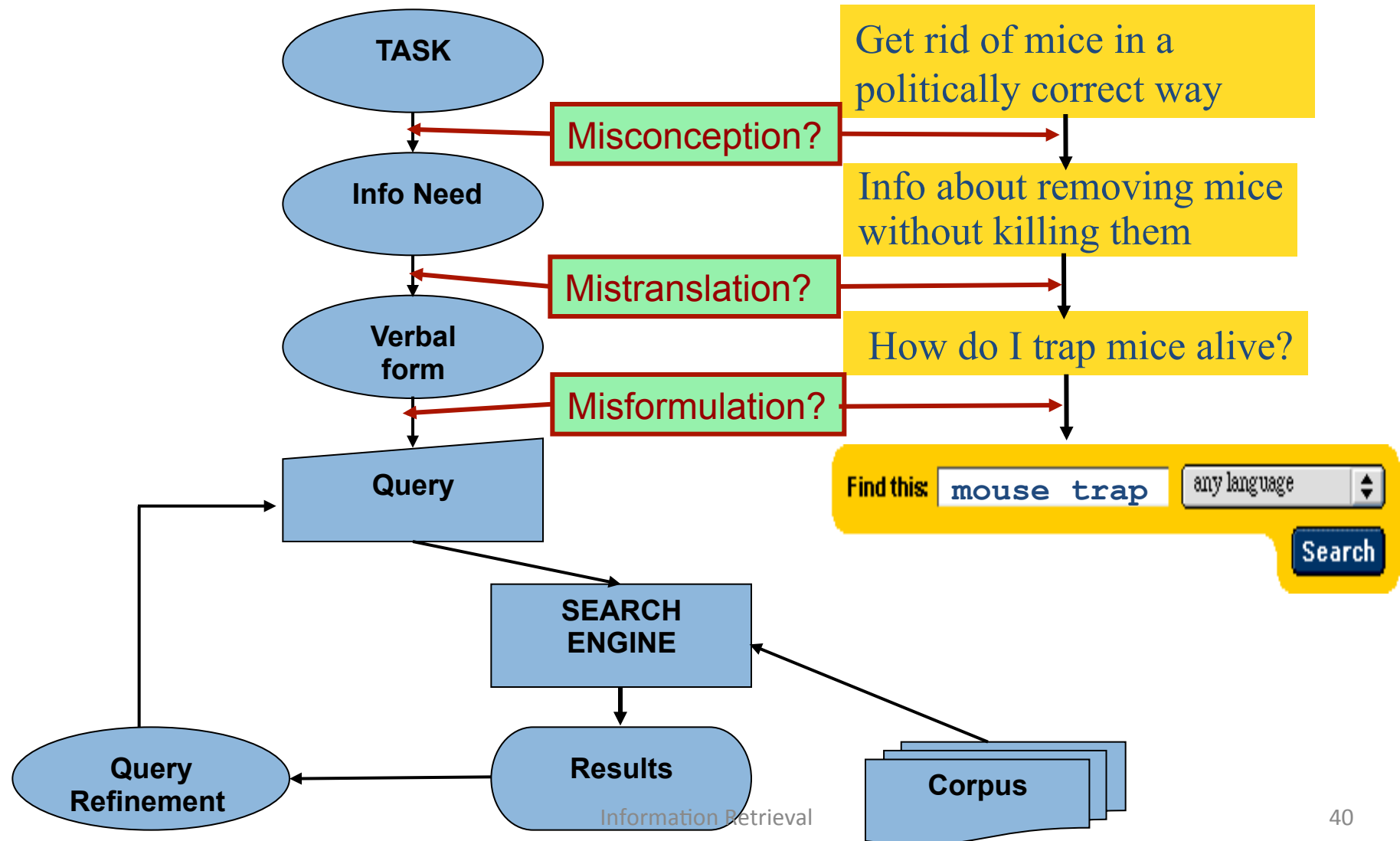
Week 13 (17 Apr) [Practical IR on the Web and Revision.](#)

1st half: about 20%

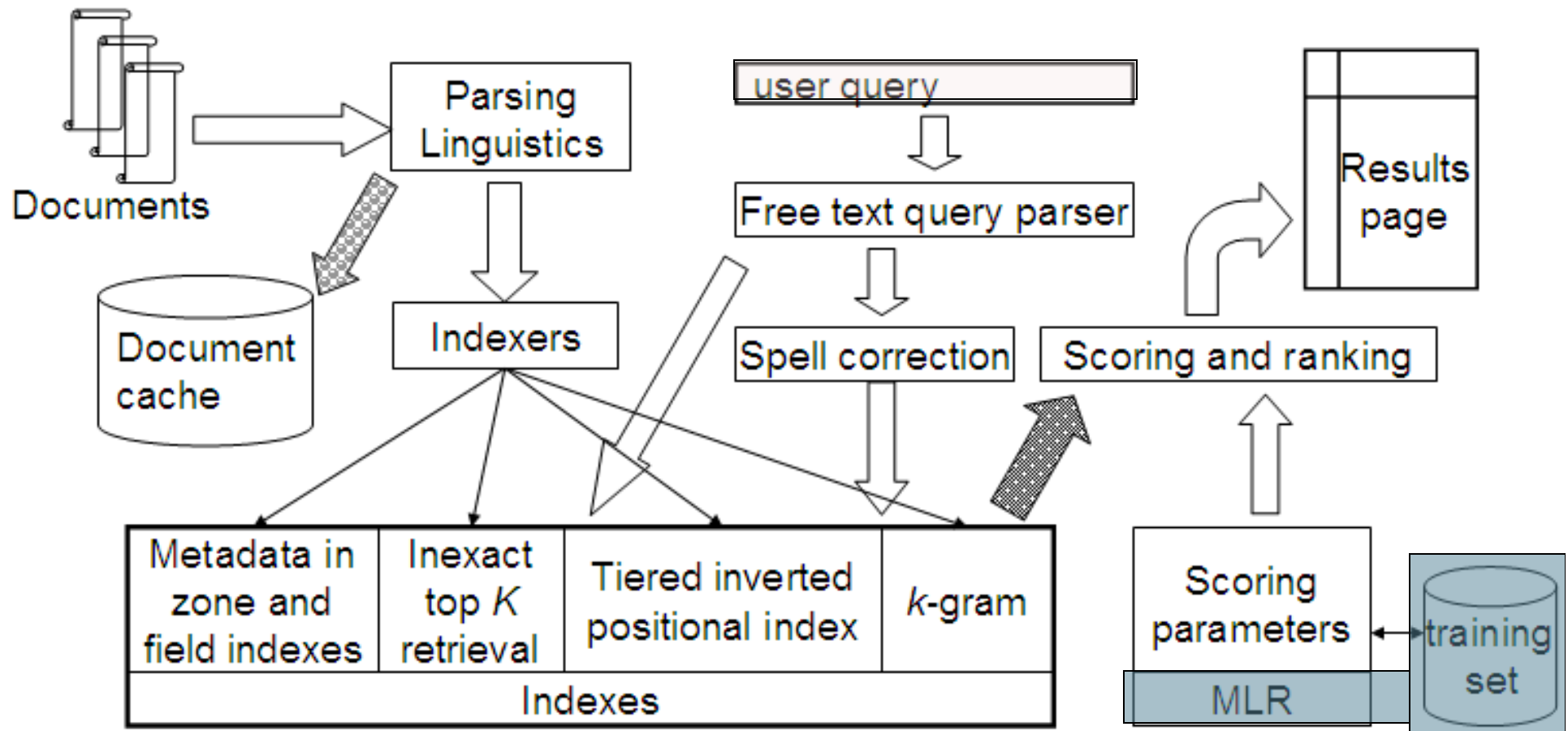
2nd half: about 80%

No homework / tutorial topics

Understanding the user: The classic search model

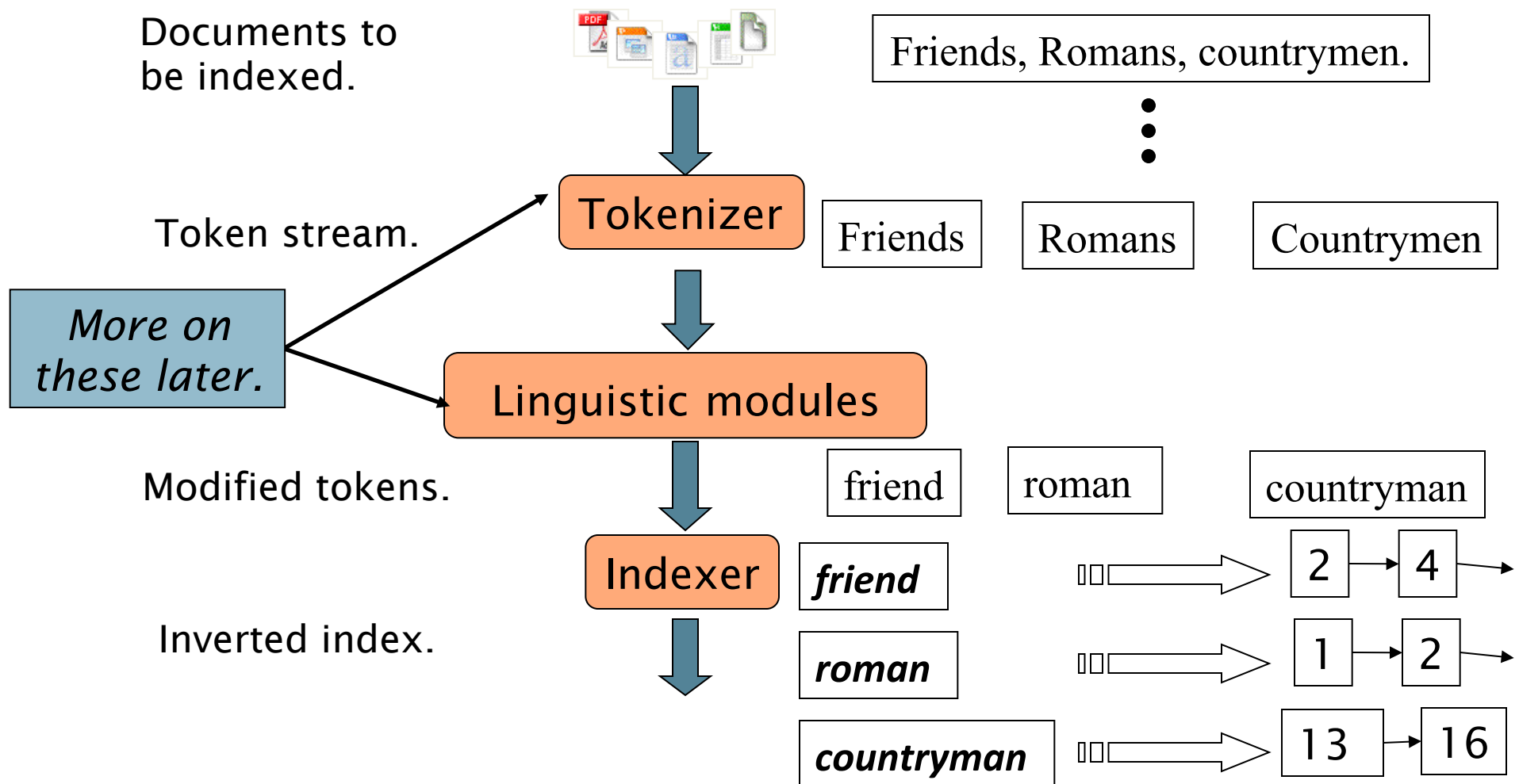


The IR System

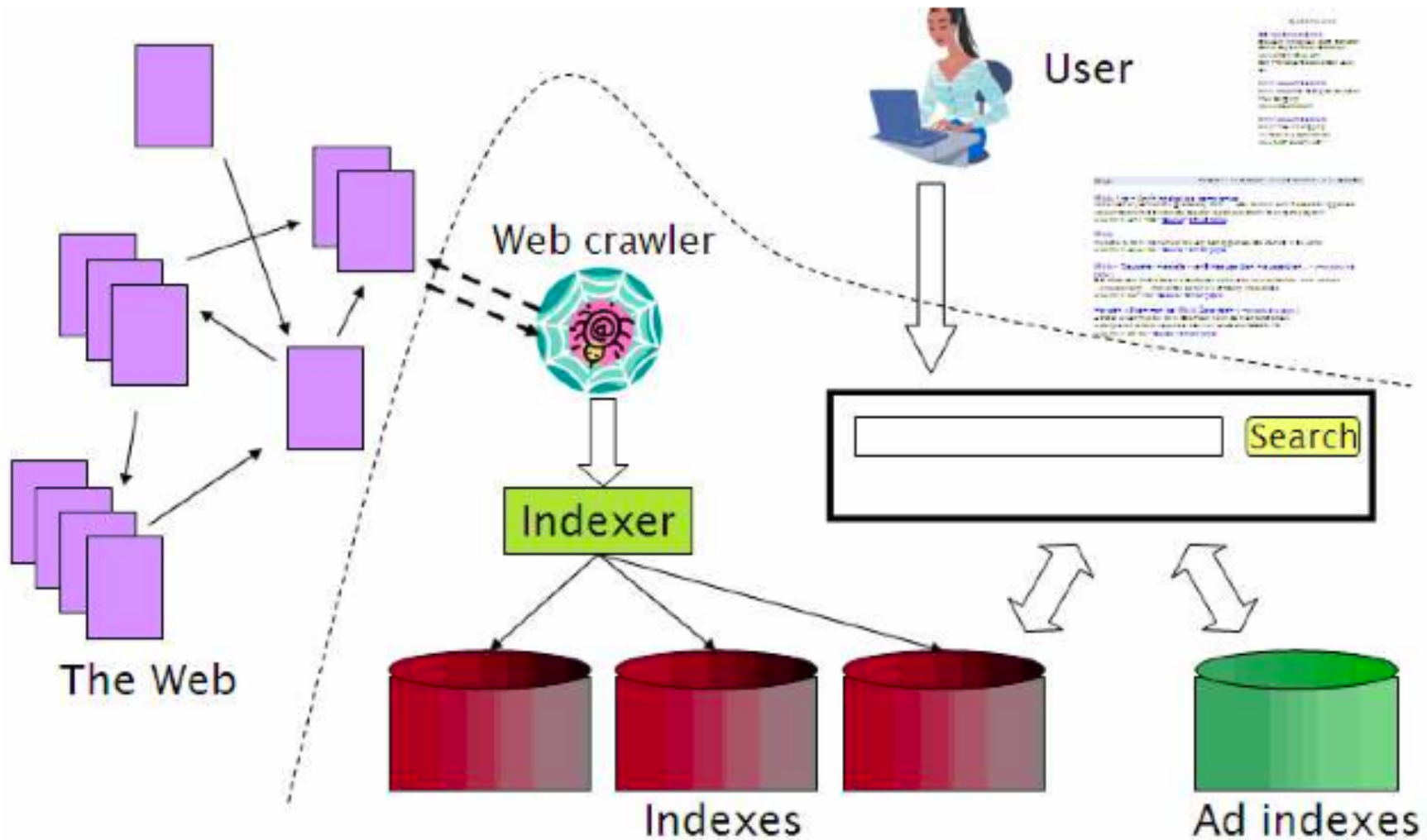


Won't be covering these blue modules in this course

Zoom in: Index Construction



Zoom Out: Web search





REVISION

Week 1: Ngram Models

- Unigram LM: Bag of words
- Ngram LM: use $n-1$ tokens of context to predict n^{th} token
- Larger n -gram models more accurate but each increase in order requires exponentially more space

Your turn: what do you think? Can we use a LM to do information retrieval?

You bet. We returned to this in Week 12.



The Unigram Model

- View language as a unordered collection of tokens
 - Each of the n tokens contributes one count (or $1/n$) to the model
 - Also known as a “bag of words”
- Outputs a count (or probability) of an input based on its individual token
 - $\text{Count}(\text{input}) = \sum_n \text{Count}(n)$
 - $P(\text{input}) = \prod_n P(n)$

Add 1 smoothing

- Not used in practice, but most basic to understand
- Idea: add 1 count to all entries in the LM, including those that are not seen

e.g., assuming $|V| = 11$

Total # of word
types in both
lines

- Q2 (By **Probability**) : “I don’t want”
 $P(\text{Aerosmith}) : .11 * .11 * .11 = 1.3\text{E-}3$
 $P(\text{LadyGaga}) : .15 * .05 * .15 = 1.1\text{E-}3$
 Winner: Aerosmith

I	2	eyes	2
don't	2	your	1
want	2	love	1
to	2	and	1
close	2	revenge	1
my	2	Total Count	18

I	3	eyes	1
don't	1	your	3
want	3	love	2
to	1	and	2
close	1	revenge	2
my	1	Total Count	20

Week 2: Basic (Boolean) IR

- Basic inverted indexes:
 - In memory dictionary and on disk postings

BRUTUS	→	1	2	4	11	31	45	173	174
--------	---	---	---	---	----	----	----	-----	-----

CAESAR	→	1	2	4	5	6	16	57	132	...
--------	---	---	---	---	---	---	----	----	-----	-----

CALPURNIA	→	2	31	54	101
-----------	---	---	----	----	-----
 - Key characteristic: Sorted order for postings
- Boolean query processing
 - Intersection by linear time “merging”
 - Simple optimizations by expected size

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

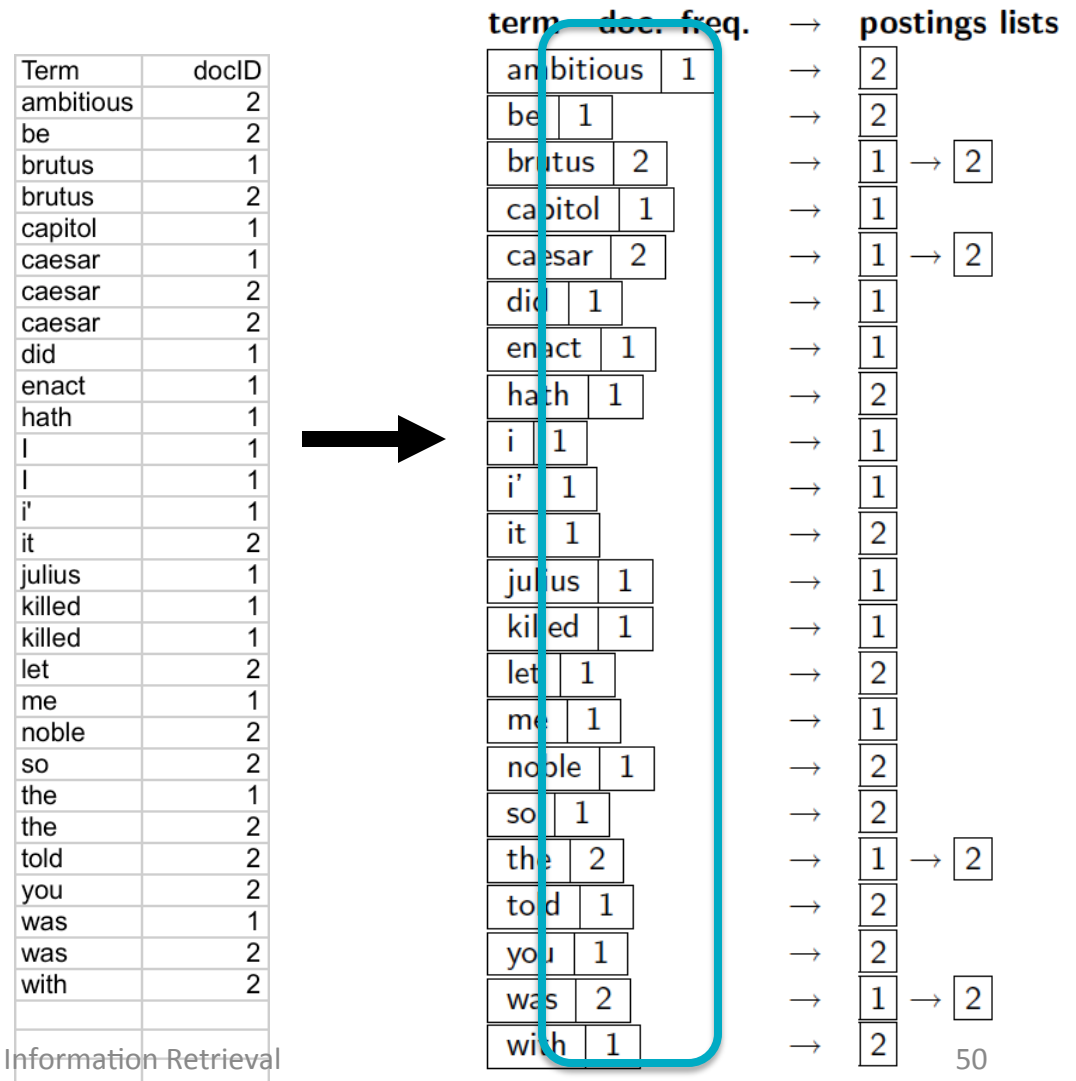
Brutus AND Caesar BUT NOT Calpurnia

1 if **play** contains **word**, 0 otherwise

Indexer steps: Dictionary & Postings

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

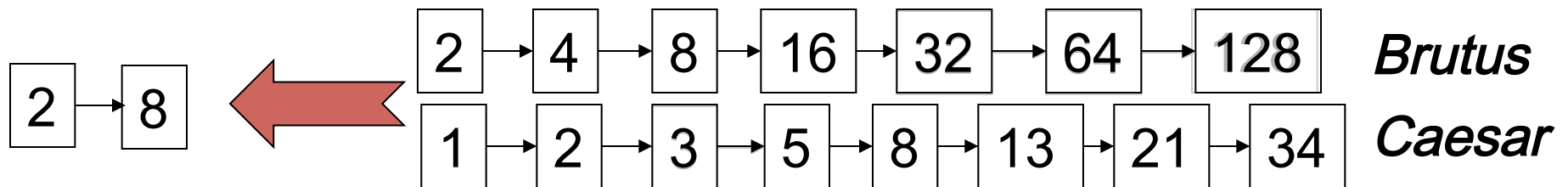
Why frequency?





The merge

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



If the list lengths are x and y , the merge takes $O(x+y)$ operations.

Crucial: postings must be sorted by docID.

Week 3: Terms and Postings Details

- The type/token distinction
 - Terms are normalized types put in the dictionary
- Tokenization problems
 - Hyphens, apostrophes, spaces, compounds
 - Language specific problems
- Term equivalence classing (or not)
 - Numbers, case folding, stemming, lemmatization
- Skip pointers
 - Encoding a tree-like structure in a postings list
- Biword indexes for phrases
- Positional indexes for phrases/proximity queries

Inverted index construction

Documents to
be indexed.



Friends, Romans, countrymen.

⋮

Token stream.

Tokenizer

Friends

Romans

Countrymen

Linguistic modules

Modified tokens.

friend

roman

countryman

Indexer

friend

→ 2 → 4 →

Inverted index.

roman

→ 1 → 2 →

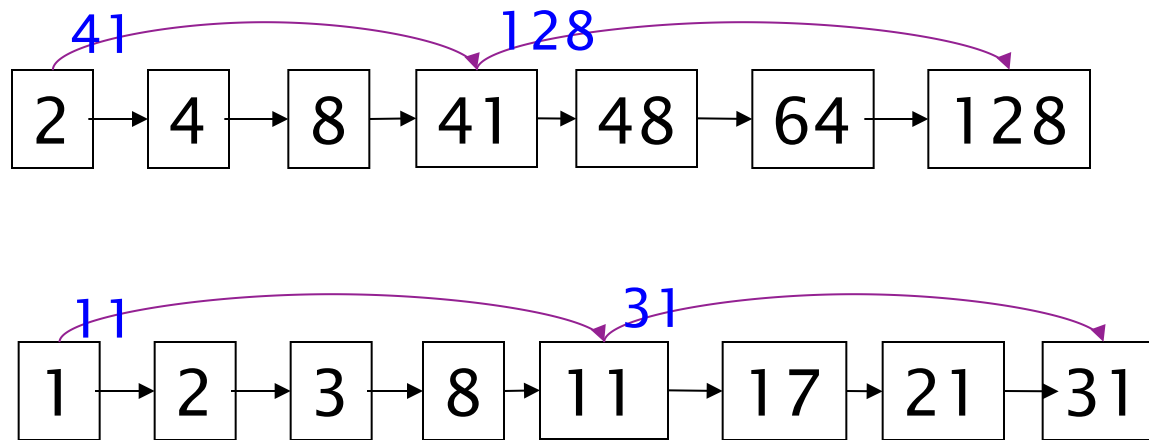
countryman

→ 13 → 16 →

Tokenization and Normalization

- Definitely language specific
- In English, we worry about
 - Tokenization – Spaces and Punctuation
 - Case folding
 - Stopwording
 - Normalization – Stemming or Lemmatization

Adding skip pointers to postings



- Done at indexing time.
- Why?
- How to do it? And where do we place skip pointers?

A first attempt at phrasal queries: Biword indexes



- Index every consecutive pair of terms in the text as a phrase: bigram model using words
- For example the text “Friends, Romans, Countrymen” would generate the biwords
 - *friends romans*
 - *romans countrymen*
- Each of these biwords is now a dictionary term
- Two-word phrase query-processing is now immediate.

Positional index example



<*be*: 993427;
1: 7, 18, 33, 72, 86, 231;
2: 3, 149;
4: 17, 191, 291, 430, 434;
5: 363, 367, ...>

Quick check:
Which of docs *1,2,4,5*
could contain "*to be*
or not to be"?

- For phrase queries, we use a merge algorithm recursively at the document level
- Now need to deal with more than just equality

Week 4: The dictionary and tolerant retrieval

- Data Structures for the Dictionary
 - Hash
 - Trees
- Learning to be tolerant
 1. Wildcards
 - General Trees
 - Permuterm
 - Ngrams, redux
 2. Spelling Correction
 - Edit Distance
 - Ngrams, re-redux
 3. Phonetic – Soundex



Hash Table

Each vocabulary term is hashed to an integer

- Pros:
 - Lookup is faster than for a tree: $O(1)$
- Cons:
 - No easy way to find minor variants:
 - judgment/judgement
 - No prefix search
 - If vocabulary keeps growing, need to occasionally do the expensive operation of rehashing *everything*

Not very tolerant!

B-trees handle *'s at the end of a query term



- How can we handle *'s in the middle of query term?
 - *co*tion*
- We could look up *co** AND **tion* in a B-tree and intersect the two term sets
 - Expensive
- The solution: transform wild-card queries so that the *'s always occur at the end
- This gives rise to the **Permuterm** Index.



Permuterm index

- For term *hello*, index under:
 - *hello\$, ello\$h, llo\$he, lo\$hel, o\$hell*
 where \$ is a special symbol.
- Queries:
 - **X** lookup on **X\$** **X*** lookup on **\$X***
 - ***X** lookup on **X\$*** ***X*** lookup on **X***
 - **X*Y** lookup on **Y\$X***

Query = hel*o
 X=hel, Y=o
 Lookup o\$hel*

Not so quick Q:
 What about X*Y*Z?



Isolated word spelling correction

- Given a lexicon and a character sequence Q , return the words in the lexicon closest to Q
- How do we define “closest”?
- We studied several alternatives
 1. Edit distance (Levenshtein distance)
 2. Weighted edit distance
 3. *n*gram overlap

Week 5: Index construction

- Sort-based indexing
 - Blocked Sort-Based Indexing
 - Merge sort is effective for disk-based sorting (avoid seeks!)
 - Single-Pass In-Memory Indexing
 - No global dictionary - Generate separate dictionary for each block
 - Don't sort postings - Accumulate postings as they occur
- Distributed indexing using MapReduce
- Dynamic indexing: Multiple indices, logarithmic merge

Hardware basics



Many design decisions in information retrieval are based on the characteristics of hardware

Especially with respect to the bottleneck:
Hard Drive Storage

- Seek Time – time to move to a random location
- Transfer Time – time to transfer a data block

BSBI: Blocked sort-based Indexing (Sorting with fewer disk seeks)

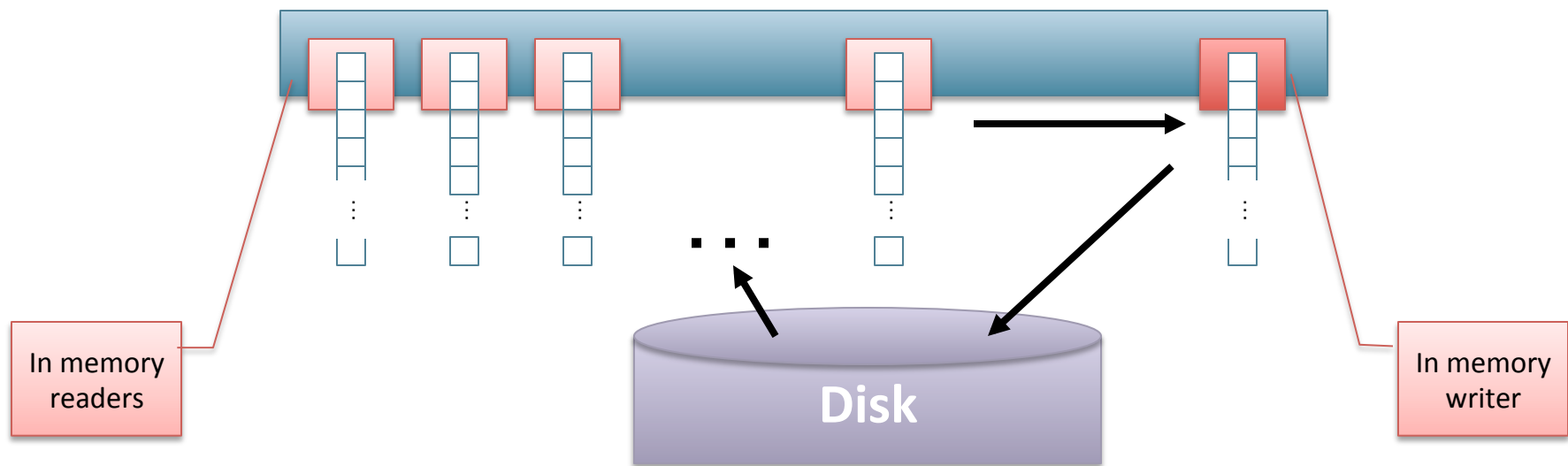


- 12-byte (4+4+4) records (*termID*, *docID*, *freq*).
- These are generated as we parse docs.
- Must now sort 100M 12-byte records by *termID*.
- Define a Block as ~ 10M such records
 - Can easily fit a couple into memory.
 - Will have 10 such blocks for our collection.
- Basic idea of algorithm:
 - Accumulate postings for each block, sort, write to disk.
 - Then merge the blocks into one long sorted order.

How to merge the sorted runs?

Second method (better):

- It is more efficient to do a n -way merge, where you are reading from all blocks simultaneously
- Providing you read decent-sized chunks of each block into memory and then write out a decent-sized output chunk, then your efficiency isn't lost by disk seeks

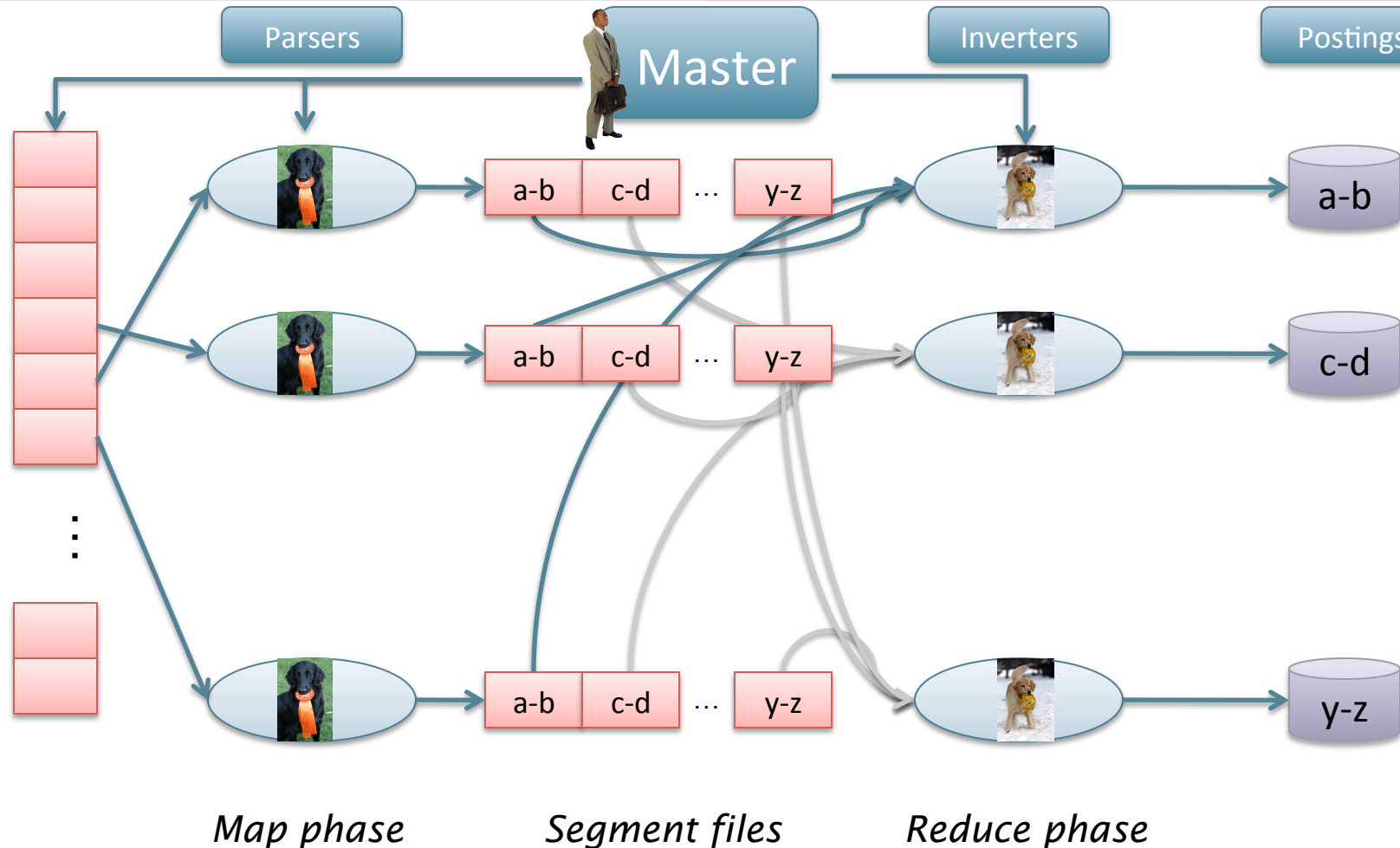


SPIMI: Single-pass in-memory indexing



- **Key idea 1:** Generate separate dictionaries for each block – no need to maintain term-termID mapping across blocks.
- **Key idea 2:** Don't sort. Accumulate postings in postings lists as they occur.
- With these two ideas we can generate a complete inverted index for each block.
- These separate indices can then be merged into one big index.

Distributed Indexing: MapReduce Data flow



Dynamic Indexing:

2nd simplest approach



- Maintain “big” main index
- New docs go into “small” (in memory) auxiliary index
- Search across both, merge results
- Deletions
 - **Invalidation bit-vector** for deleted docs
 - Filter docs output on a search result by this invalidation bit-vector
- Periodically, re-index into one main index
 - Assuming T total # of postings and n as size of auxiliary index, we touch each posting **up to** $\text{floor}(T/n)$ times.



Logarithmic merge

- Idea: maintain a series of indexes, each twice as large as the previous one.
- Keep smallest (Z_0) in memory
- Larger ones (I_0, I_1, \dots) on disk
- If Z_0 gets too big ($> n$), write to disk as I_0 or merge with I_0 (if I_0 already exists) as Z_1
- Either write merge Z_1 to disk as I_1 (if no I_1) Or merge with I_1 to form Z_2
- ... etc.

Loop for log levels

Week 6: Index Compression

- Collection and vocabulary statistics: Heaps' and Zipf's laws

Compression to make index smaller, faster

- Dictionary compression for Boolean indexes
 - Dictionary string, blocks, front coding
- Postings compression: Gap encoding

Empirical Laws



Heaps' law: $M = kT^b$

- M is the size of the vocabulary, T is the number of tokens in the collection
- In a log-log plot of vocabulary size M vs. T , Heaps' law predicts a line with slope about $\frac{1}{2}$
 - It is the simplest possible relationship between the two in log-log space

Zipf's law: $cf_i \propto 1/i = K/i$

- Zipf's law: The i th most frequent term has frequency proportional to $1/i$.
- where K is a normalizing constant
- cf_i is collection frequency (not document frequency): the number of occurrences of the term t_i in the collection.

Index Compression: Dictionary-as-a-String and Blocking



- Store pointers to every k th term string.
 - Example below: $k=4$.
- Need to store term lengths (1 extra byte)

....**7***systile***9***syzygetic***8***syzygial***6***syzygy***11***szaibelyite* ...

Freq.	Postings ptr.	Term ptr.
33		
29		
44		
126		
7		

} Save 9 bytes
on 3
pointers.

← Lose 4 bytes on
term lengths.

Postings Compression: Postings file entry



- We store the list of docs containing a term in increasing order of docID.
 - **computer**: 33,47,154,159,202 ...
- Consequence: it suffices to store *gaps*.
 - 33,14,107,5,43 ...
- Hope: most gaps can be encoded/stored with far fewer than 20 bits.

Variable Byte Encoding: Example

docIDs	824	829	215406
gaps		5	214577
VB code	00000110 10111000	10000101	00001101 00001100 10110001

$$512 + 256 + 32 + 16 + 8 = 824$$

Postings stored as the byte concatenation

00000110 10111000 10000101 00001101 00001100 10110001

Key property: VB-encoded postings are uniquely prefix-decodable.

For a small gap (5), VB uses a whole byte.

Week 7: Vector space ranking

1. Represent the query as a weighted tf-idf vector
2. Represent each document as a weighted tf-idf vector
3. Compute the cosine similarity score for the query vector and each document vector
4. Rank documents with respect to the query by score
5. Return the top K (e.g., $K = 10$) to the user

Term-document count matrices

- Store 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

tf-idf weighting



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

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

- **Best known weighting scheme IR**
 - 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



Queries as vectors

- Key idea 1: Do the same for queries: represent them as vectors in the space; they are “mini-documents”
- Key idea 2: Rank documents according to their proximity to the query in this space
- proximity = similarity of vectors
- proximity \approx inverse of distance
- **Motivation: Want to get away from the you’re-either-in-or-out Boolean model.**
- Instead: rank more relevant documents higher than less relevant documents



Length normalization

- A vector can be (length-) normalized by dividing each of its components by its length – for this we use the L_2 norm:

$$\|\vec{x}\|_2 = \sqrt{\sum_i x_i^2}$$

- Dividing a vector by its L_2 norm makes it a unit (length) vector (on surface of unit hypersphere)
- Effect on the two documents d and d' (d appended to itself) from earlier slide: they have identical vectors after length normalization.
 - Long and short documents now have comparable weights

Cosine for length-normalized vectors

- For length-normalized vectors, cosine similarity is simply the dot product (or scalar product):

$$\cos(\vec{q}, \vec{d}) = \vec{q} \cdot \vec{d} = \sum_{i=1}^{|V|} q_i d_i$$

for q, d length-normalized.

Week 8: Complete Search Systems

Making the Vector Space Model more efficient to compute

- Approximating the actual correct results
- Skipping unnecessary documents

In actual data: dealing with zones and fields, query term proximity

Resources for today

- IIR 7, 6.1

Recap: Computing cosine scores

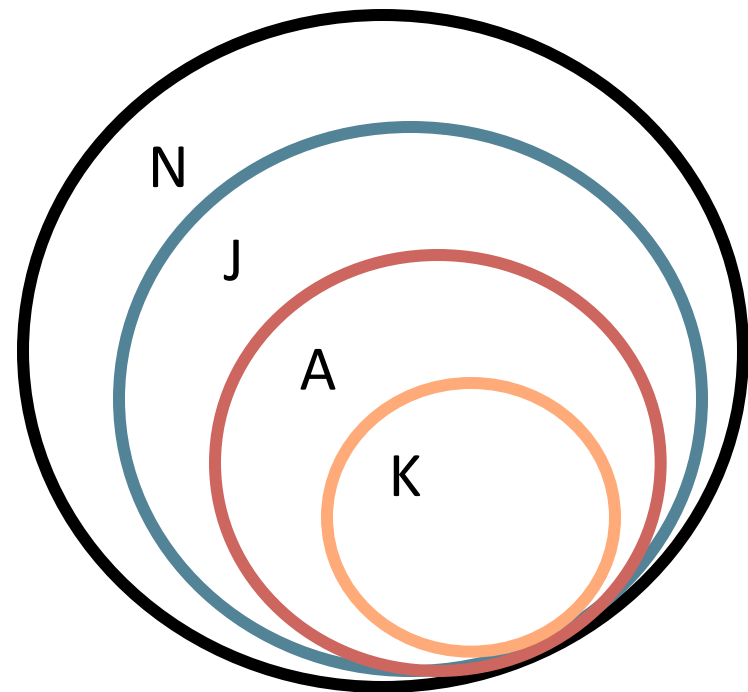
COSINESCORE(q)

```
1  float Scores[N] = 0
2  float Length[N]
3  for each query term  $t$ 
4  do calculate  $w_{t,q}$  and fetch postings list for  $t$ 
5      for each pair( $d, tf_{t,d}$ ) in postings list
6      do  $Scores[d] += w_{t,d} \times w_{t,q}$ 
7  Read the array  $Length$ 
8  for each  $d$ 
9  do  $Scores[d] = Scores[d] / Length[d]$ 
10 return Top  $K$  components of  $Scores[]$ 
```



Generic approach

- Find a set A of *contenders*, with $K < |A| \ll N$
 - A does not necessarily contain the top K , but has many docs from among the top K
 - Return the top K docs in A
- Think of A as pruning non-contenders
- The same approach can also be used for other (non-cosine) scoring functions





Net score

- Consider a simple total score combining cosine relevance and authority

$$\text{net-score}(q,d) = g(d) + \text{cosine}(q,d)$$

- Can use some other linear combination than an equal weighting
 - Indeed, any function of the two “signals” of user happiness
-
- Now we seek the top K docs by net score



Parametric Indices

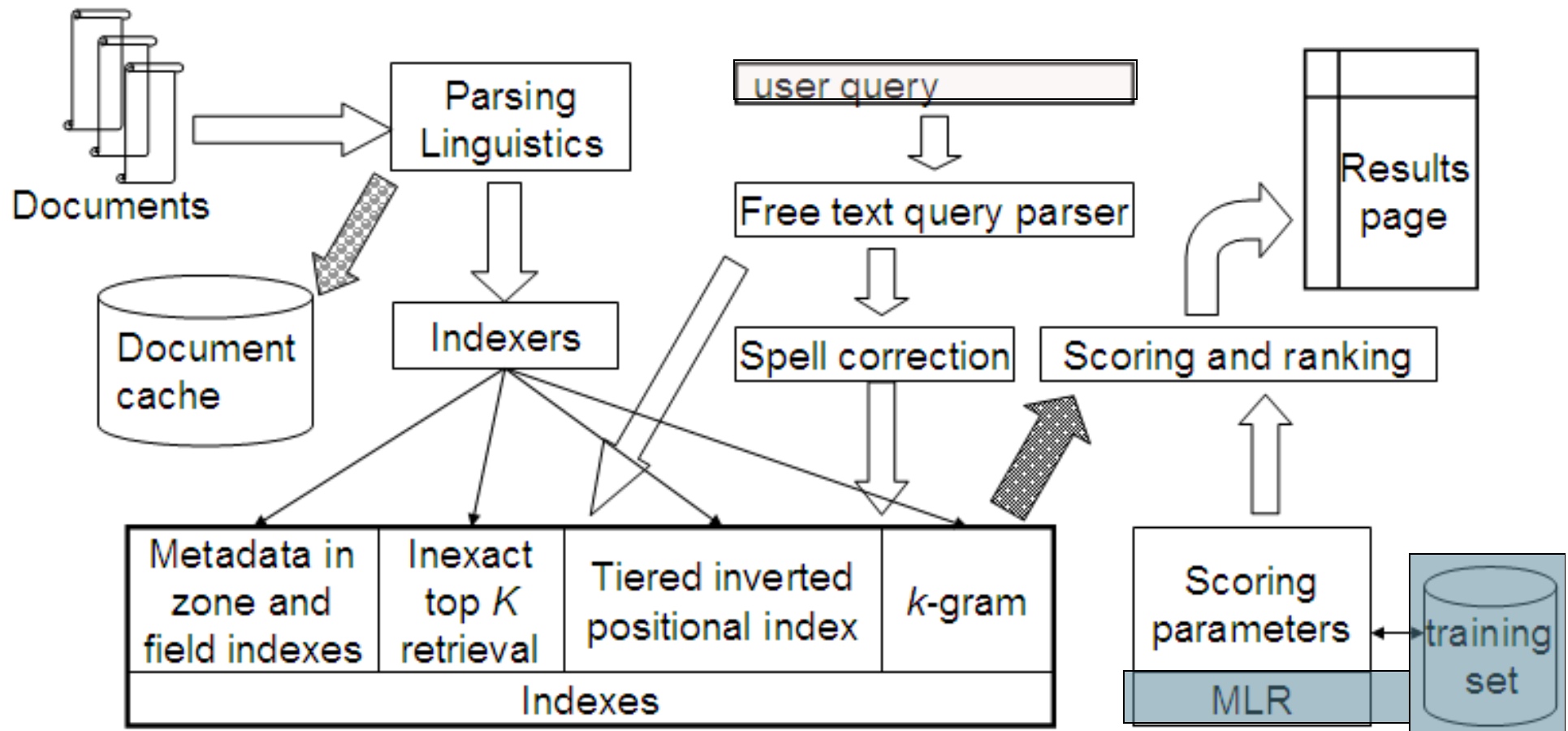
Fields

- Year = 1601 is an example of a field
- Field or parametric index: postings for each field value
 - Sometimes build range (B-tree) trees (e.g., for dates)
- Field query typically treated as conjunction
 - (doc *must* be authored by shakespeare)

Zone

- A zone is a region of the doc that can contain an arbitrary amount of text
e.g.,
 - Title
 - Abstract
 - References ...
- Build inverted indexes on zones as well to permit querying

Putting it all together

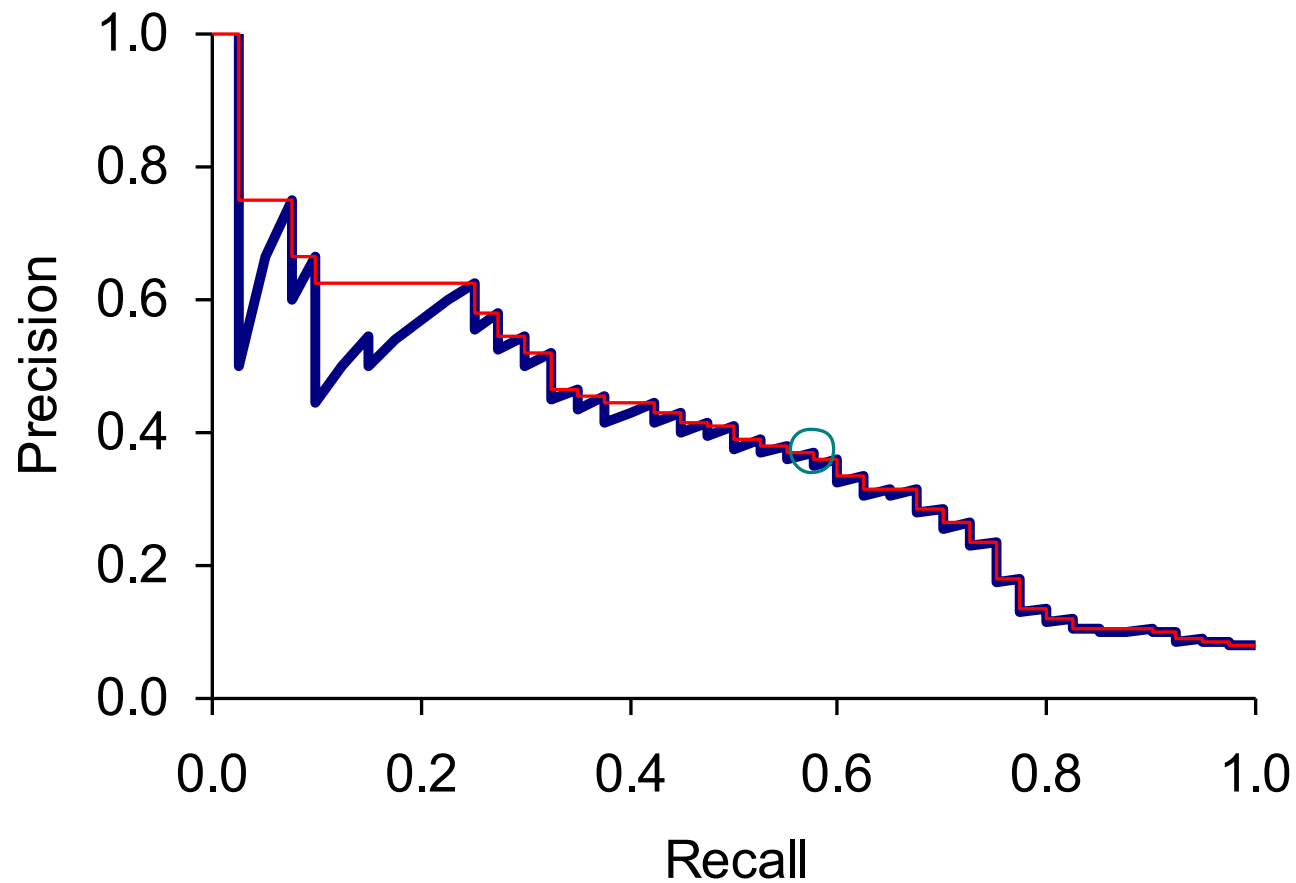


Won't be covering these blue modules in this course

Week 9: IR Evaluation

- How do we know if our results are any good?
 - Evaluating a search engine
 - Benchmarks
 - Precision and Recall; Composite measures

A precision-recall curve





Kappa Example

$$P(A) = 370/400 = 0.925$$

$$P(\text{nonrelevant}) = (10+20+70+70)/800 = 0.2125$$

$$P(\text{relevant}) = (10+20+300+300)/800 = 0.7878$$

$$P(E) = 0.2125^2 + 0.7878^2 = 0.665$$

$$\text{Kappa} = (0.925 - 0.665)/(1-0.665) = 0.776$$

- $\text{Kappa} > 0.8 \rightarrow$ good agreement
- $0.67 < \text{Kappa} < 0.8 \rightarrow$ “tentative conclusions”
- Depends on purpose of study
- For >2 judges: average pairwise kappas

A/B testing



Purpose: Test a single innovation




Prerequisite: You have a large system with many users.

- Have most users use old system
- Divert a small proportion of traffic (e.g., 1%) to the new system that includes the innovation
- Evaluate with an “automatic” overall evaluation criterion (OEC) like clickthrough on first result
- Now we can directly see if the innovation does improve user happiness.
- Probably the evaluation methodology that large search engines trust most



Dynamic summaries

- Present one or more “windows” within the document that contain several of the query terms
 - One of the killer features of Google (ca. 1996)
 - “KWIC” snippets: Keyword in Context presentation

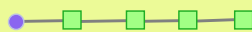
 <input type="text" value="christopher manning"/>	<u>Christopher Manning, Stanford NLP</u> Christopher Manning , Associate Professor of Computer Science and Linguistics, Stanford University. nlp.stanford.edu/~manning/ - 12k - Cached - Similar pages
 <input type="text" value="christopher manning machine translation"/>	<u>Christopher Manning, Stanford NLP</u> Christopher Manning , Associate Professor of Computer Science and Linguistics, ... computational semantics, machine translation , grammar induction, ... nlp.stanford.edu/~manning/ - 12k - Cached - Similar pages
 <input type="text" value="christopher manning"/>	<u>Christopher Manning, Stanford NLP</u> Christopher Manning , Associate Professor of Computer Science and Linguistics, Stanford University ... Chris Manning works on systems and formalisms that can ... nlp.stanford.edu/~manning/ - Cached

Kinds of behaviors we see in the data

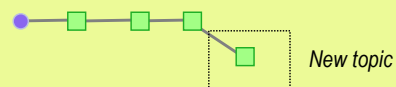
Short / Nav



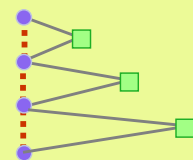
Topic exploration



Topic switch



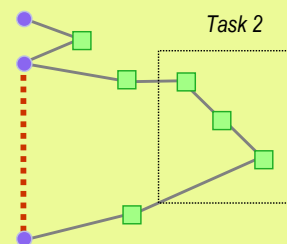
Methodical results
exploration



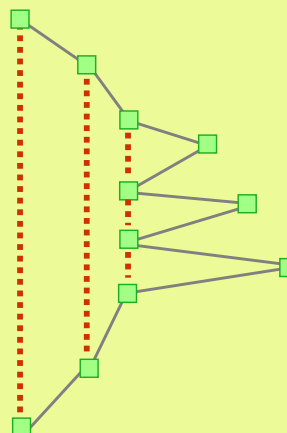
Query reform



Multitasking



Stacking behavior



38

Week 10: XML, Relevance Feedback and Query Expansion

Chapter 10

1. XML

- Basic XML concepts
- Challenges in XML IR
- Vector space model for XML IR

Chapter 9

1. Relevance Feedback

Document Level

- Explicit RF – Rocchio (1971)
- When does it work?
- Variants – Implicit and Blind

2. Query Expansion

Term Level

- Controlled Vocabularies
- WordNet
- Automatic Thesaurus Generation

XML Basics and Definitions

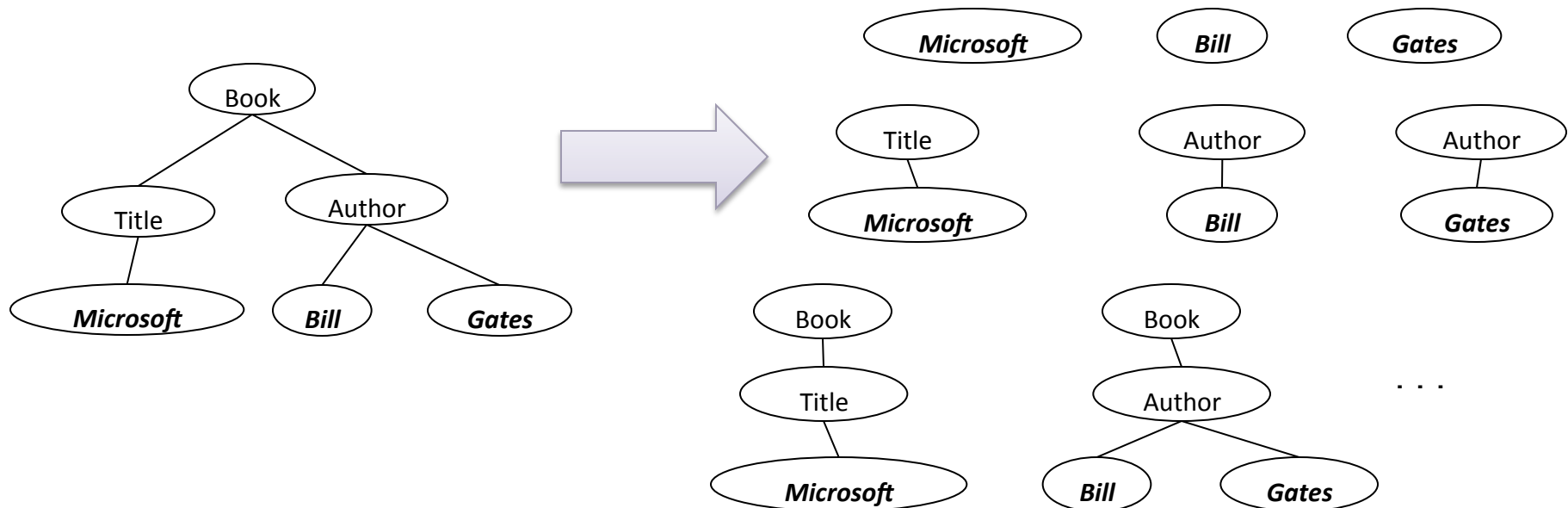


- **XML Document Object Model (XML DOM):** standard for accessing and processing XML documents
 - The DOM represents elements, attributes and text within elements as nodes in a tree.
 - With a DOM API, we can process an XML documents by starting at the root element and then descending down the tree from parents to children.
- **XPath:** standard for enumerating path in an XML document collection.
 - We will also refer to paths as **XML contexts** or simply **contexts**
- **Schema:** puts constraints on the structure of allowable XML documents. E.g. a schema for Shakespeare's plays: scenes can occur as children of acts.
 - Two standards for schemas for XML documents are: XML DTD (document type definition) and XML Schema.

Main idea: lexicalized subtrees

- Aim: to have each dimension of the vector space encode a word together with its position within the XML tree.
- How: Map XML documents to lexicalized subtrees.

“With words”





Context resemblance

- A simple measure of the similarity of a path c_q in a query and a path c_d in a document is the following **context resemblance** function CR:

$$\text{CR}(c_q, c_d) = \begin{cases} \frac{1+|c_q|}{1+|c_d|} & \text{if } c_q \text{ matches } c_d \\ 0 & \text{if } c_q \text{ does not match } c_d \end{cases}$$

$|c_q|$ and $|c_d|$ are the number of nodes in the query path and document path, respectively

- c_q matches c_d **iff** we can transform c_q into c_d by inserting additional nodes.

INEX relevance assessments

- The relevance-coverage combinations are quantized as follows:

$$Q(rel, cov) = \begin{cases} 1.00 & \text{if } (rel, cov) = 3E \\ 0.75 & \text{if } (rel, cov) \in \{2E, 3L\} \\ 0.50 & \text{if } (rel, cov) \in \{1E, 2L, 2S\} \\ 0.25 & \text{if } (rel, cov) \in \{1S, 1L\} \\ 0.00 & \text{if } (rel, cov) = 0N \end{cases}$$

- This evaluation scheme takes account of the fact that binary relevance judgments are not appropriate for XML retrieval. The quantization function **Q** instead allows us to grade each component as partially relevant. The number of relevant components in a retrieved set *A* of components can then be computed as:

$$\#(\text{relevant items retrieved}) = \sum_{c \in A} Q(rel(c), cov(c))$$


Rocchio (1971)



Popularized in the SMART system (Salton)

In practice:

$$\vec{q}_m = \alpha \vec{q}_0 + \beta \frac{1}{|D_r|} \sum_{\vec{d}_j \in D_r} \vec{d}_j - \gamma \frac{1}{|D_{nr}|} \sum_{\vec{d}_j \in D_{nr}} \vec{d}_j$$

- D_r = set of known relevant doc vectors
- D_{nr} = set of known irrelevant doc vectors
-  Different from C_r and C_{nr} as we only get judgments from a few documents
- $\{\alpha, \beta, \gamma\}$ = weights (hand-chosen or set empirically)

Query Expansion



- In relevance feedback, users give additional input (relevant/non-relevant) on **documents**, which is used to reweight terms in the documents
- In query expansion, users give additional input (good/bad search term) on **words or phrases**

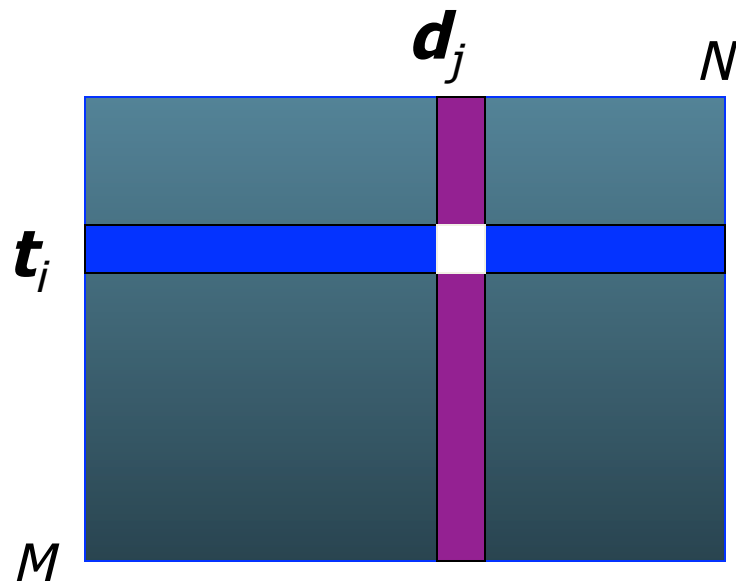


Co-occurrence Thesaurus



Simplest way to compute one is based on term-term similarities in $C = AA^T$ where A is term-document matrix.

- $w_{i,j}$ = (normalized) weight for (t_i, d_j)



- For each t_i , pick terms with high values in C

In NLTK. Did you forget?

A concordance permits us to see words in context. For example, we saw that then inserting the relevant word in parentheses:

```
>>> text1.similar("monstrous")
Building word-context index...
subtly impalpable pitiable curious imperial perilous trust
abundant untoward singular lamentable few maddens horrible
mystifying christian exasperate puzzled
>>> text2.similar("monstrous")
Building word-context index...
very exceedingly so heartily a great good amazingly as swee
remarkably extremely vast
>>>
```

Observe that we get different results for different texts. Austen uses this word

The term `common_contexts` allows us to examine just the contexts that are sh

```
>>> text2.common_contexts(["monstrous", "very"])
be_glad am_glad a_pretty is_pretty a_lucky
>>>
```

Week 11: Probabilistic IR

Chapter 11

1. Probabilistic Approach to Retrieval / Basic Probability Theory
2. Probability Ranking Principle
3. OKAPI BM25

Chapter 12

1. Language Models for IR



Binary Independence Model (BIM)

- Traditionally used with the PRP

Assumptions:

- **Binary** (equivalent to Boolean): documents and queries represented as binary term incidence vectors
 - E.g., document d represented by vector $\vec{x} = (x_1, \dots, x_M)$, where
 - $x_t = 1$ if term t occurs in d and $x_t = 0$ otherwise
 - Different documents may have the same vector representation
- **Independence**: no association between terms (not true, but works in practice – **naïve** assumption)



Okapi BM25: A Nonbinary Model

- If the query is long, we might also use similar weighting for **query terms**

$$RSV_d = \sum_{t \in q} \left[\log \frac{N}{df_t} \right] \cdot \frac{(k_1 + 1)tf_{td}}{k_1((1 - b) + b \times (L_d/L_{ave})) + tf_{td}} \cdot \frac{(k_3 + 1)tf_{tq}}{k_3 + tf_{tq}}$$

- tf_{tq} : term frequency in the query q
- k_3 : tuning parameter controlling term frequency scaling of the query
- No length normalization of queries (because retrieval is being done with respect to a single fixed query)
- The above tuning parameters should be set by optimization on a development test collection. Experiments have shown reasonable values for k_1 and k_3 as values between 1.2 and 2 and $b = 0.75$



An Appraisal of Probabilistic Models

- The difference between ‘vector space’ and ‘probabilistic’ IR is not that great:
 - In either case you build an information retrieval scheme in the exact same way.
 - Difference: for probabilistic IR, at the end, you score queries not by cosine similarity and tf-idf in a vector space, but by a slightly different formula motivated by probability theory

Using language models in IR



- Each document is treated as (the basis for) a language model.
- Given a query q , rank documents based on $P(d|q)$
- $P(q)$ is the same for all documents, so ignore
- $P(d)$ is the prior – often treated as the same for all d
 - But we can give a prior to “high-quality” documents, e.g., those with high static quality score $g(d)$ (cf. Section 7.14).
- $P(q|d)$ is the probability of q given d .
- So to rank documents according to relevance to q , ranking according to $P(q|d)$ and $P(d|q)$ is equivalent.

Mixture model: Summary



$$P(q|d) \propto \prod_{1 \leq k \leq |q|} (\lambda P(t_k|M_d) + (1 - \lambda)P(t_k|M_c))$$

- What we model: The user has a document in mind and generates the query from this document.
- The equation represents the probability that the document that the user had in mind was in fact this one.

Week 12: Web Search

Chapter 20

- Crawling

Chapter 21

- Anchor Text
- PageRank



IR on the web vs. IR in general

- On the web, search is not just a nice feature.
 - Search is a key enabler of the web: financing, content creation, interest aggregation, etc.
- look at search ads
- The web is a chaotic und uncoordinated collection.
 - lots of duplicates – need to detect duplicates
- No control / restrictions on who can author content.
 - lots of spam – need to detect spam
- The web is very large. → need to know how big it is.



Pagerank summary

- Pre-processing:
 - Given graph of links, build matrix **A**
 - From it compute **a**
 - The pagerank a_i is a scaled number between 0 and 1
- Query processing:
 - Retrieve pages meeting query
 - Rank them by their pagerank
 - Order is *query-independent*



PageRank issues

- Real surfers are not random surfers.
 - Examples of nonrandom surfing: back button, short vs. long paths, bookmarks, directories – and search!
 - → Markov model is not a good model of surfing.
 - But it's good enough as a model for our purposes.
- Simple PageRank ranking (as described on previous slide) produces bad results for many pages.
 - Consider the query [video service].
 - The Yahoo home page (i) has a very high PageRank and (ii) contains both *video* and *service*.
 - If we rank all Boolean hits according to PageRank, then the Yahoo home page would be top-ranked.
 - Clearly not desirable.



**WHERE TO GO
FROM HERE**



Learning Objectives

In addition to learning about IR, you have picked up skills that you will help in your future computing

- **Python** – one of the easiest and more straightforward programming languages to use.
- **NLTK** – A good set of routines and data that are useful in dealing with NLP and IR.



Computational
Advertising

Geographic IR

Adversarial IR

Distributed IR

Query
Analysis

Social
Network IR

IR Theory

Natural Language
Processing

Digital Libraries

Question
Answering

Prob IR

Recommendation
Systems

VSM

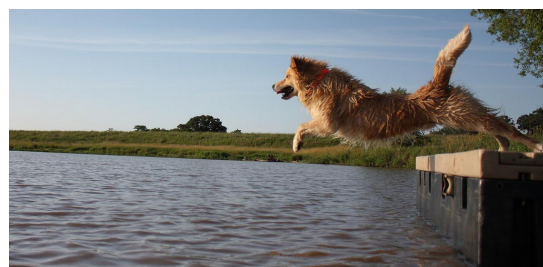
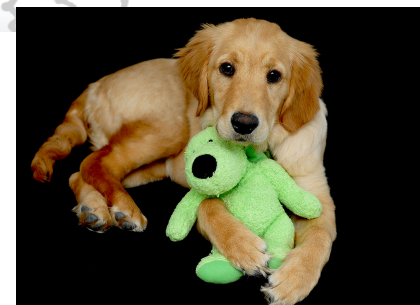
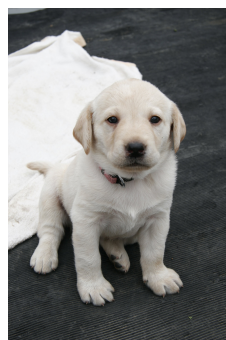
Boolean IR

Photo credits:
<http://www.flickr.com/photos/aperezdc/>



Opportunities in IR





Thanks for joining us for the journey!

