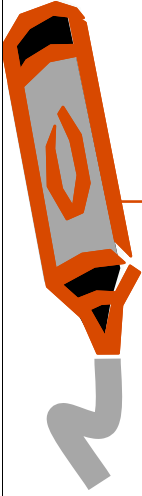




Information Retrieval

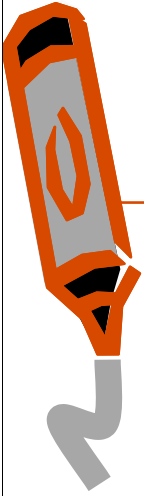
Overview

many slides courtesy James Allan @umass Amherst
some slides courtesy ChengXiang Zhai @Urbana-Champaign



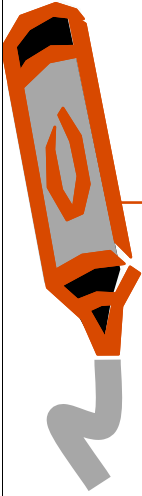
administrivia

- INSTRUCTOR
Virgil Pavlu vip@ccs.neu.edu
- GRADER : if we get any, he/she will have final authority on homework grading issues
- <http://ccs.neu.edu/~jaa/ISU535.06X2>



overview

- what is information retrieval ?
- how does it work ?
- a simple retrieval model



Google™

Web Images Groups News Froogle Local Scholar more »

ikim 2003

Search

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

what is IR?

Web

[12. CIKM 2003: New Orleans, Louisiana, USA](#)

12. **CIKM 2003**: New Orleans, Louisiana, USA. Proceedings of the **2003 ACM CIKM** International Conference on Information and Knowledge Management, New Orleans, ...
www.informatik.uni-trier.de/~ley/db/conf/cikm/cikm2003.html - 56k - [Cached](#) - [Similar pages](#)

[CIKM](#)

Proceedings of the **2003 ACM CIKM** International Conference on Information and Knowledge Management, New Orleans, Louisiana, USA, November 2-8, **2003**. ...
www.informatik.uni-trier.de/~ley/db/conf/cikm/ - 10k - Jun 25, 2005 - [Cached](#) - [Similar pages](#)

[CIKM'2003 review](#)

CIKM'2003 highlights. 12th ACM International Conference on Information and Knowledge Management, 3-8 November, New Orleans ...
smi.ucd.ie/~rinat/papers/cikm03_rep.html - 22k - [Cached](#) - [Similar pages](#)

[Collaborative Filtering Mailing List Archive: \[collab@sims\] CFP](#)

ACM **CIKM 2003** Call For Papers. 12th International Conference on Information and Knowledge ... caliber papers submitted to **CIKM 2003** will be accepted. ...
www.pdesigner.net/1996/0697.html - 17k - [Cached](#) - [Similar pages](#)

[TOC](#)

Proceedings of the twelfth international conference on Information and knowledge management citation. **2003**, New Orleans, LA, USA November 03 - 08, **2003** ...
portal.acm.org/toc.cfm?id=956863&type=proceeding - [Similar pages](#)

[\[Asis-\] CIKM 2003](#)

[Asis-] **CIKM 2003**. Padmini Srinivasan padmini@lakshmi.info-science.uiowa.edu Mon, 29 Sep 2003 12:59:36 -0500. Previous message: [Asis-] Re: ...
mail.asis.org/pipermail/asis-l/2003-September/001024.html - 17k - [Cached](#) - [Similar pages](#)

[\[PDF\] CIKM 2003](#)

File Format: PDF/Adobe Acrobat - [View as HTML](#)
CIKM 2003. Jacob Kogan. Charles Nicholas. Marc Teboulle. -means and beyond - p.1/53. Page 2. Outline of the talk. how to build a partition ...
www.csee.umbc.edu/~nicholas/clustering/jacob.pdf - [Similar pages](#)

[Tutorial on Document Clustering](#)

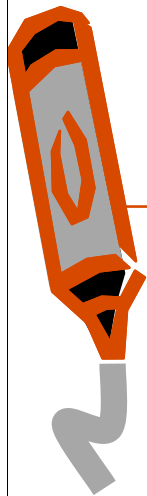
CIKM 2003 Tutorial. Clustering Large and High-Dimensional Data ... Katya Pelekhev and Daniela Rus,"Using Star Clusters for Filtering", **CIKM 2000**, (pdf) ...
www.csee.umbc.edu/~nicholas/clustering/ - 9k - [Cached](#) - [Similar pages](#)

[Conference on Information and Knowledge Management \(CIKM\)](#)

CIKM has a strong tradition of workshops devoted to emerging areas of database ... The **CIKM 2004** web page; The **CIKM 2003** Web Page; The **CIKM 2002** Web Page ...
www.cikm.org/ - 7k - [Cached](#) - [Similar pages](#)

[CIKM 2003, New Orleans, USA, November 2003](#)

Home. **CIKM 2003**, New Orleans, USA, November 2003. << Bild 6 | Bild 7/80 | Bild 8 >>. Miniaturansicht.
www.torsten-priebe.de/showpics.php?folder=2003-11a_cikm03&picture=7 - 2k - [Cached](#) - [Similar pages](#)



Google

Web Images Groups News Froogle Local Scholar more »

cikm 2003

Search

Advanced Search Preferences

evaluation

Web

[12. CIKM 2003: New Orleans, Louisiana, USA](#)

12. **CIKM 2003**: New Orleans, Louisiana, USA. Proceedings of the **2003 ACM CIKM** International Conference on Information and Knowledge Management, New Orleans, ...
www.informatik.uni-trier.de/~ley/db/conf/cikm/cikm2003.html - 56k - [Cached](#) - [Similar pages](#)



[CIKM](#)

Proceedings of the **2003 ACM CIKM** International Conference on Information and Knowledge Management, New Orleans, Louisiana, USA, November 2-8, **2003**. ...
www.informatik.uni-trier.de/~ley/db/conf/cikm/ - 10k - Jun 25, 2005 - [Cached](#) - [Similar pages](#)



[CIKM'2003 review](#)

CIKM'2003 highlights. 12th ACM International Conference on Information and Knowledge Management, 3-8 November, New Orleans ...
smi.ucd.ie/~rinat/papers/cikm03_rep.html - 22k - [Cached](#) - [Similar pages](#)



[Collaborative Filtering Mailing List Archive: \[collab@sims\] CFP](#)

ACM **CIKM 2003** Call For Papers. 12th International Conference on Information and Knowledge ... caliber papers submitted to **CIKM 2003** will be accepted. ...
www.pdesigner.net/1996/0697.html - 17k - [Cached](#) - [Similar pages](#)



[TOC](#)

Proceedings of the twelfth international conference on Information and knowledge management citation. **2003**, New Orleans, LA, USA November 03 - 08, **2003** ...
portal.acm.org/toc.cfm?id=956863&type=proceeding - [Similar pages](#)



[\[Asis-l\] CIKM 2003](#)

[Asis-l] **CIKM 2003**. Padmini Srinivasan padmini@lakshmi.info-science.uiowa.edu Mon, 29 Sep 2003 12:59:36 -0500. Previous message: [Asis-l] Re: ...
mail.asis.org/pipermail/asis-l/2003-September/001024.html - 17k - [Cached](#) - [Similar pages](#)



[\[PDF\] CIKM 2003](#)

File Format: PDF/Adobe Acrobat - [View as HTML](#)
CIKM 2003. Jacob Kogan. Charles Nicholas. Marc Tebouille. -means and beyond - p.1/53. Page 2. Outline of the talk. how to build a partition ...
www.csee.umbc.edu/~nicholas/clustering/jacob.pdf - [Similar pages](#)



[Tutorial on Document Clustering](#)

CIKM 2003 Tutorial. Clustering Large and High-Dimensional Data ... Katya Pelekhov and Daniela Rus,"Using Star Clusters for Filtering", **CIKM** 2000, (pdf) ...
www.csee.umbc.edu/~nicholas/clustering/ - 9k - [Cached](#) - [Similar pages](#)



[Conference on Information and Knowledge Management \(CIKM\)](#)

CIKM has a strong tradition of workshops devoted to emerging areas of database ... The **CIKM 2004** web page; The **CIKM 2003** Web Page; The **CIKM 2002** Web Page ...
www.cikm.org/ - 7k - [Cached](#) - [Similar pages](#)

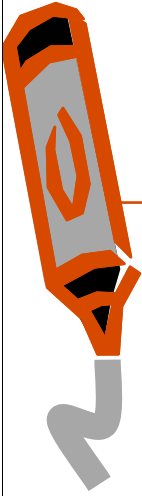


[CIKM 2003, New Orleans, USA, November 2003](#)

Home. **CIKM 2003**, New Orleans, USA, November 2003. << Bild 6 | Bild 7/80 | Bild 8 >>. Miniaturansicht.



www.torsten-priebe.de/showpics.php?folder=2003-11a_cikm03&picture=7 - 2k - [Cached](#) - [Similar pages](#)



altavista Web Images MP3/Audio Video News

cikm 2003 [Advanced Search](#)
[Settings](#)

SEARCH: Worldwide USA RESULTS IN: All languages English, Spanish

web search

AltaVista found 19,700 results

[Conference on Information and Knowledge Management \(CIKM\)](#)

... The **CIKM 2004** web page. The **CIKM 2003** Web Page. The **CIKM 2002** Web Page. The **CIKM 2001** Web page
www.cikm.org
[More pages from cikm.org](#)

[Department of Computer Science at LSU](#)

News & Events. Department of Computer Science. 298 Coates Hall. Phone: (225)578-1495
bit.csc.lsu.edu
[More pages from bit.csc.lsu.edu](#)

[ACM CIKM 2003 Call For Papers](#)

ACM **CIKM 2003** –Call for Industry Track Presentations. We solicit high-quality Industry Track presentati
management products and marketplace trends ... Web site: <http://www.cikm.org/2003>. Sponsored by ACM SIG
www.cs.wisc.edu/dbworld/messages/2003-06/1057176548.html
[More pages from cs.wisc.edu](#)

[MMDB'03](#)

... conjunction with The 12th International Conference on Information and Knowledge Management (ACM **CIK**
Call for ...
www.cs.fiu.edu/mmdb03
[More pages from cs.fiu.edu](#)

[CIKM 2002 Homepage](#)

The **CIKM 2002** Technical Program has been updated. Hotel information has been updated in the Local Arran
international forum for presentation ... The 12th International Conference on Information and Knowledge Mana
cikm.org/2002
[More pages from cikm.org](#)

[DBWorld Message](#)

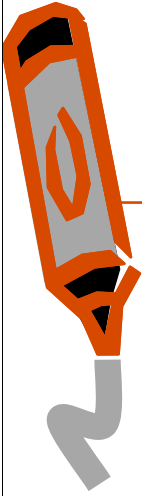
Only the highest caliber papers submitted to **CIKM 2003** will be accepted. We have a special interest in papers
For Papers 12th International Conference on Information and Knowledge Management (**CIKM'03** ...
www.cs.wisc.edu/dbworld/messages/2003-03/1051827399.html
[More pages from cs.wisc.edu](#)

[Past Events](#)

... July 28 - August 1, **2003**. **SIGIR 2003** (Toronto) - 26th Annual International ACM SIGIR Conference ... May 28,
www.sigir.org/events/events-past.html
[More pages from sigir.org](#)

[Upcoming Conference Deadlines](#)

... 5/20/**2003**. ACM **CIKM 2003**. Electronic abstracts deadline [Website] ... 5/28/**2003**. **CIKM 2003**. Paper submis
www-users.cs.umn.edu/~gladmin/conferences
[More pages from www-users.cs.umn.edu](#)



alltheweb advanced search :: customize preferences :: submit site :: help
cikm 2003 SEARCH
• • • find it all • • • Results in: Any Language English
Web News Pictures Video Audio
1 - 10 of 18,400 Results for cikm 2003

web search

Web Results [\(What's this?\)](#)

[Conference on Information and Knowledge Management \(CIKM\)](#)

... The **CIKM 2004** web page. The **CIKM 2003** Web Page. The **CIKM 2002** Web Page. The **CIKM 2001** Web page. TI
[more hits from:](#) <http://www.cikm.org/> - 6 KB

[12. CIKM 2003: New Orleans, Louisiana, USA](#)

CIKM 2003: New Orleans, Louisiana, USA. Proceedings of the **2003 ACM CIKM** International Conference on Inform
ISBN 1-58113-723-0

[more hits from:](#) <http://www.informatik.uni-trier.de/~ley/db/conf/cikm/cikm2003.html> - 55 KB

[Department of Computer Science at LSU](#)

News & Events. Department of Computer Science. 298 Coates Hall. Phone: (225)578-1495

[more hits from:](#) <http://bit.csc.lsu.edu/> - 9 KB

[CIKM](#)

International Conference on Information and Knowledge Management (**CIKM**) 14. **CIKM 2005:** Bremen, Germany. ·
CIKM 2003 Home Page. 11. **CIKM 2002:** McLean, Virginia, USA ...

[more hits from:](#) <http://sunsite.informatik.rwth-aachen.de/dblp/db/conf/cikm> - 10 KB

[ACM WIDM'2003](#)

CIKM 2003. Past WIDM's. Important Dates: Submission of abstract. June 28, **2003**. Submission full paper: July 5,
Past WIDM's ... Conference on Information and Knowledge Management (**CIKM 2003**) Hotel Inter-Continental, Nev

[more hits from:](#) <http://www.cais.ntu.edu.sg/widm2003> - 12 KB

[ACM CIKM 2003 Call For Papers](#)

ACM **CIKM 2003** Call for Industry Track Presentations. We solicit high-quality Industry Track presentations which p
marketplace trends ... Web site: <http://www.cikm.org/2003>. Sponsored by ACM SIGIR and ACM SIGMIS (pending a

[more hits from:](#) <http://www.cs.wisc.edu/dbworld/messages/2003-06/1057176548.html> - 39 KB

[ACM CIKM 2003 :: Agents Portal :: Agent and Multi-agent Technology Resources](#)

... PRELIMINARY CALL FOR PAPERS. <http://cikm.org/2003/> Twelfth International Conference on Information and ...

[more hits from:](#) <http://aose.ift.ulaval.ca/modules.php?op=modload&name=News&file=article&sid=86> - 27 KB

[MMDB'03](#)

... conjunction with The 12th International Conference on Information and Knowledge Management (ACM **CIKM 2003**)

[more hits from:](#) <http://www.cs.fiu.edu/mmdb03> - 6 KB

[CIKM 2002 Homepage](#)

The **CIKM 2002** Technical Program has been updated. Hotel information has been updated in the Local Arrangem
for presentation ... The 12th International Conference on Information and Knowledge Management (**CIKM 2003**) ...

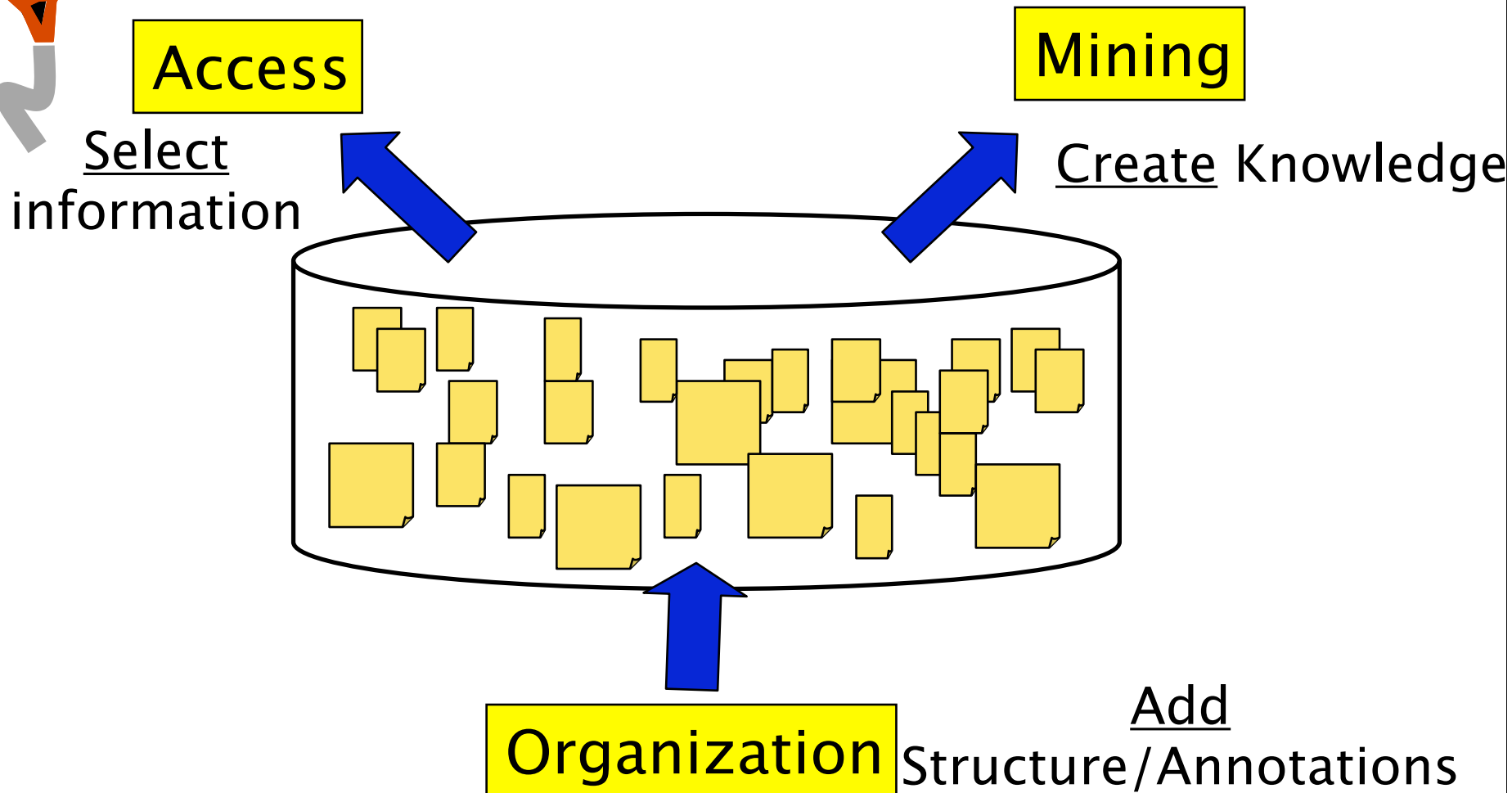
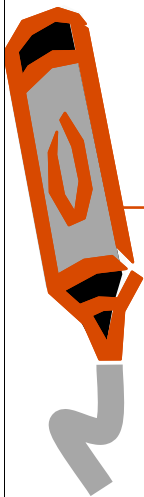
[more hits from:](#) <http://cikm.org/2002> - 15 KB

[DBWorld Message](#)

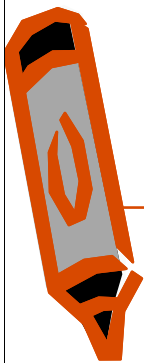
Only the highest caliber papers submitted to **CIKM 2003** will be accepted. We have a special interest in papers the
International Conference on Information and Knowledge Management (**CIKM'03**) ...

[more hits from:](#) <http://www.cs.wisc.edu/dbworld/messages/2003-03/1051827399.html> - 4 KB

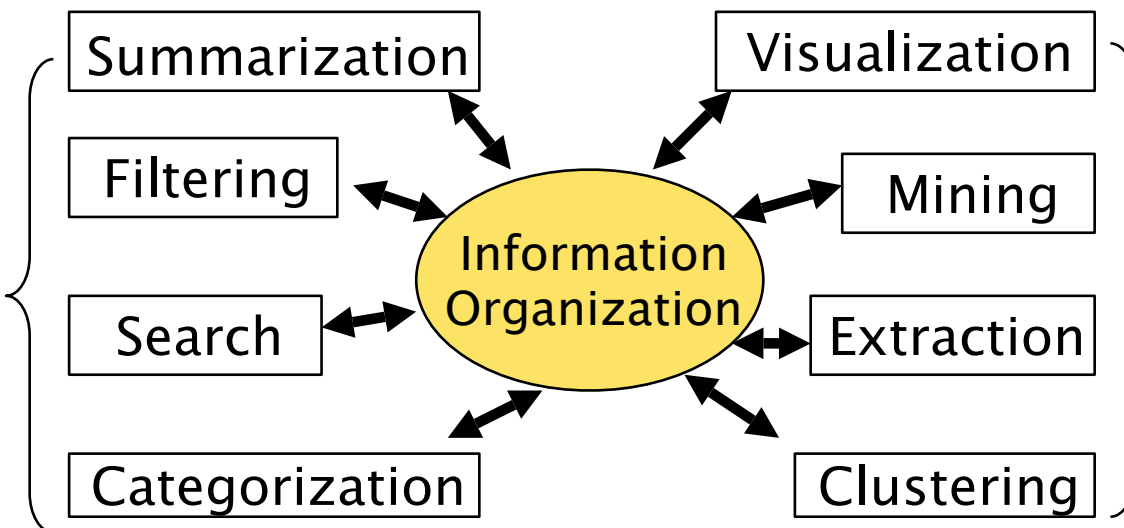
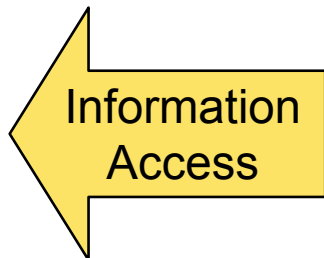
text management applications



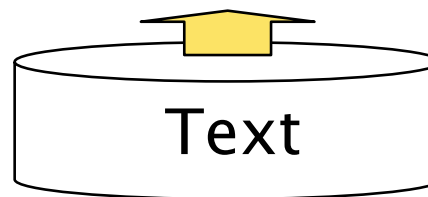
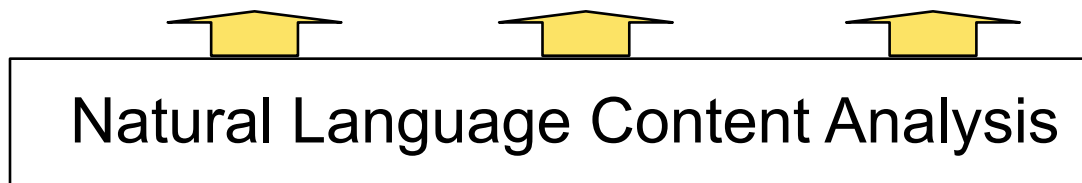
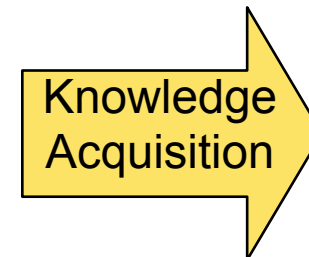
text management applications



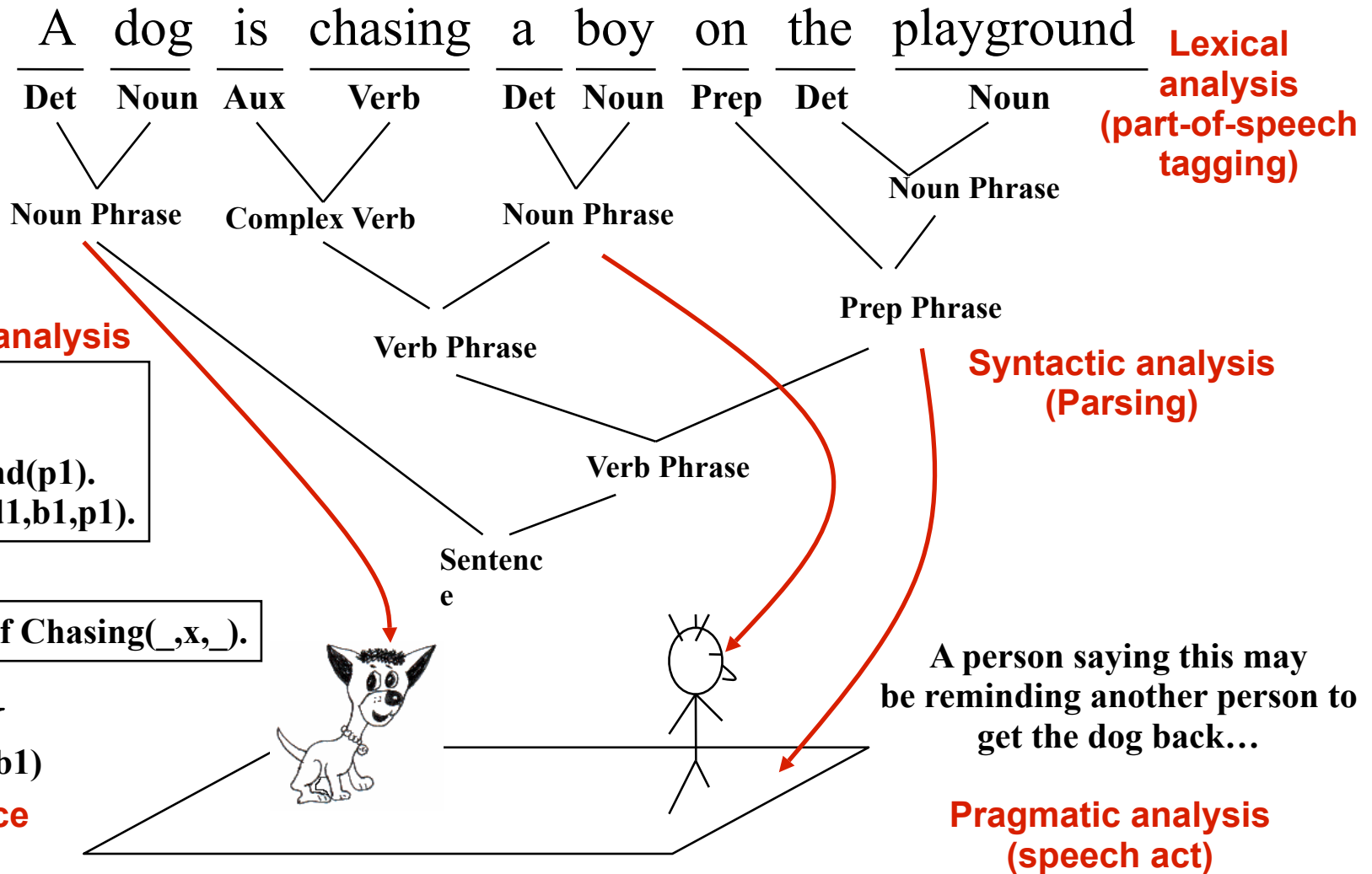
Retrieval Applications



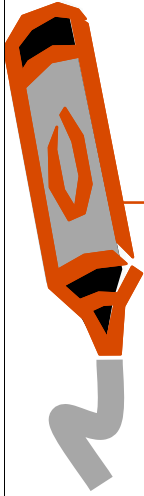
Mining Applications



natural language processing



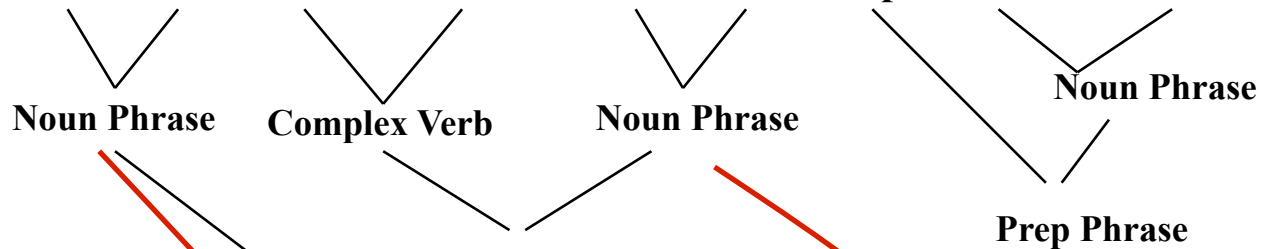
what we can do in NLP



A dog is chasing a boy on the playground

Det Noun Aux Verb Det Noun Prep Det Noun

POS Tagging: 97%



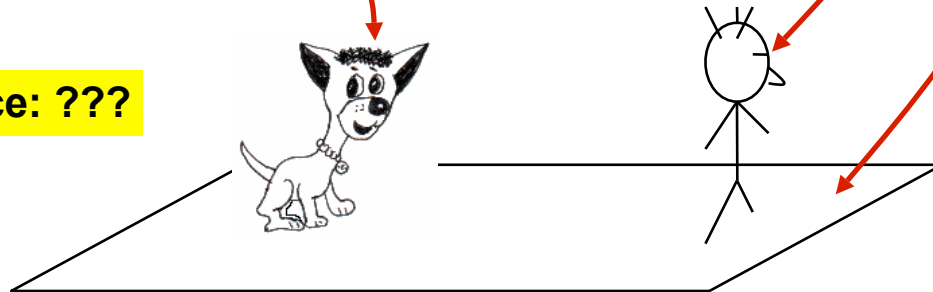
Parsing: partial >90%(?)

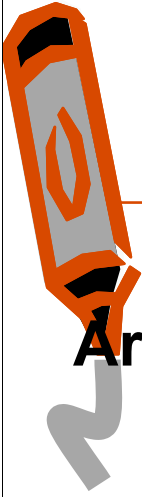
Semantics: some aspects

- Entity/relation extraction
- Word sense disambiguation
- Anaphora resolution

Speech act analysis: ???

Inference: ???





natural language processing

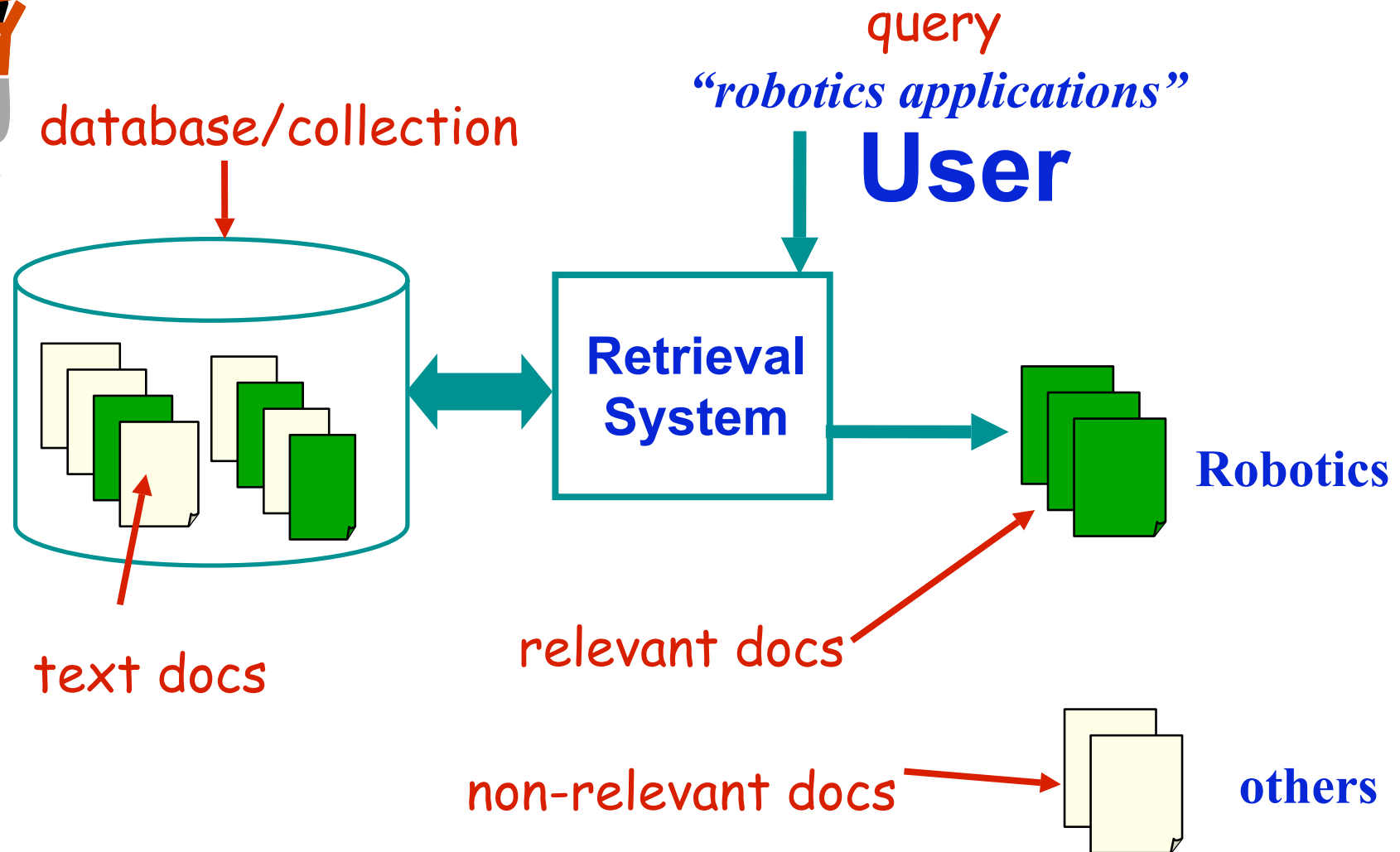
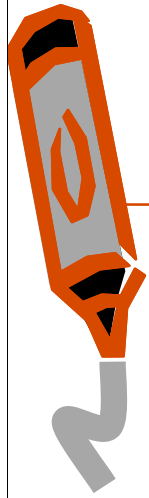
Arabic text

... يَجِبُ عَلَى الْإِنْسَانِ أَنْ يَكُونَ أَمِينًا وَصَادِقًا مَعَهُ
نَفْسِهِ وَمَعَ أَهْلِهِ وَجَارِيَّتِهِ وَأَنْ يَبْذُلَ كُلَّ جُودٍ فِي إِعْلَاءِ
شَأْنِ الْوَطَنِ وَأَنْ يَعْمَلَ عَلَى مَآ...

How can a computer make **sense** out of this **string**?

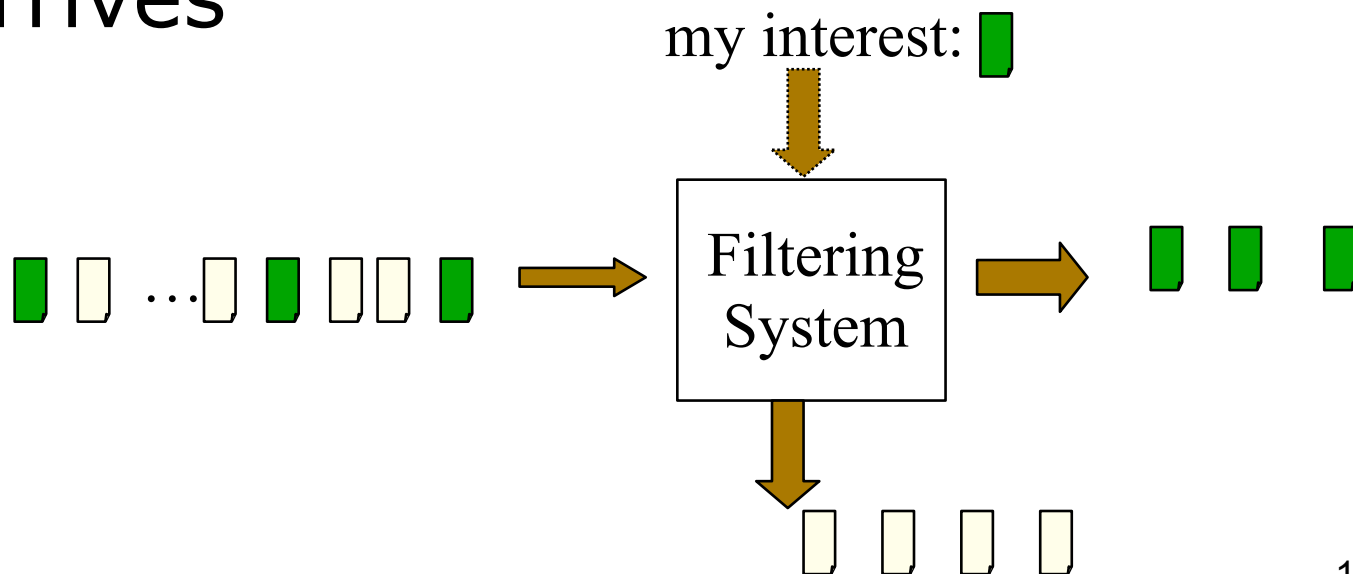
- Morphology** - What are the basic units of meaning (words)?
- What is the meaning of each word?
- Syntax** - How are words related with each other?
- Semantics** - What is the “combined meaning” of words?
- Pragmatics** - What is the “meta-meaning”? (speech act)
- Discourse** - Handling a large chunk of text
- Inference** - Making sense of everything

search (ad-hoc IR)

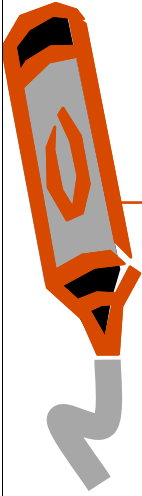


information filtering

- Stable & long term interest, dynamic info source
- System must make a delivery decision immediately as a document "arrives"



collaborative filtering



amazon.com Your Store Books See All 31 Product Categories Your Account | Cart | Wish List | Help


Search Books [GO] Browse Subjects | Bestsellers | The New York Times® Best Sellers | Magazines | Corporate Accounts | E-books & Docs | Bargain Books | Used Books

Advanced Search

Your order qualifies for free shipping! (Some restrictions apply)

Make sure to select **FREE Super Saver Shipping** as your shipping speed at checkout.

You could save \$30 today with the Amazon Visa® Card:

 Your current subtotal: \$109.76
Amazon Visa discount: - \$30.00
Your new subtotal: **\$79.76** [Find out how](#)

Save \$30 off your first purchase, earn 3% rewards, get a 0% APR*, and pay no annual fee.

Customers who bought *Managing Gigabytes* also bought:



[Mining the Web](#)
by Soumen Chakrabarti

Price: **\$57.95**
Used & new from \$41.42

[Add to cart](#)

[Explore similar items](#)



[Foundations of Statistical Natural Language Processing](#)
by Christopher D. Manning, Hinrich Schtze

Price: **\$67.32**
Used & new from \$42.98

[Add to cart](#)



[Natural Language Processing for Online Applications](#)
by Peter Jackson, Isabelle Moulinier

Price: **\$126.00**
Used & new from \$142.44

[Add to cart](#)

Customers who shopped for *Managing Gigabytes* also shopped for:



[Information Retrieval](#)
by William B. Frakes, Ricardo Baeza-Yates

Price: **\$69.67**
Used & new from \$37.03

[Add to cart](#)

[Explore similar items](#)



[Understanding Search Engines](#)
by Michael W. Berry, Murray Browne

Price: **\$41.50**
Used & new from \$28.50

[Add to cart](#)



[Survey of Text Mining](#)
by Michael W. Berry (Editor)

Price: **\$58.93**
Used & new from \$54.59

[Add to cart](#)

Want free shipping? You're almost there! Add a recommended item and qualify now. [Some restrictions apply.](#)



[Natural Language Processing for Online Applications](#)
by Peter Jackson, Isabelle Moulinier

Price: **\$39.95**
Used & new from \$44.89



[Lucene in Action \(In Action series\)](#)
by Erik Hatcher, Otis Gospodnetic

Price: **\$29.67**
Used & new from \$28.00

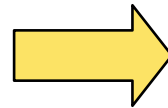
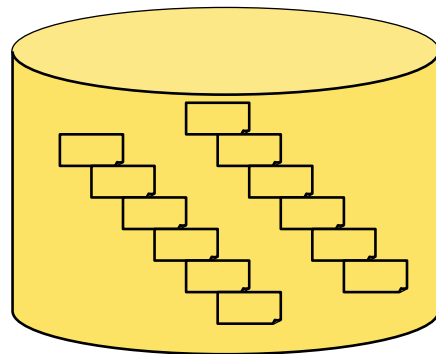


[Speech and Language Processing](#)
by Daniel Jurafsky, James H. Martin

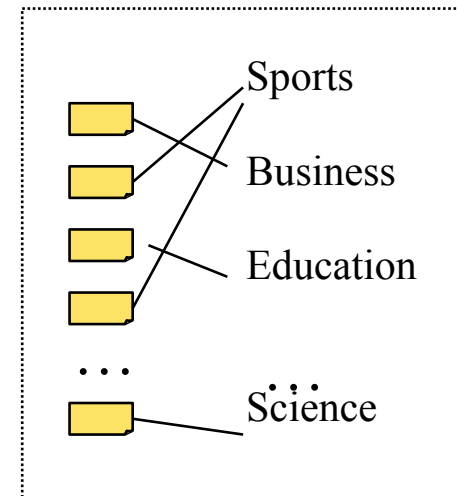
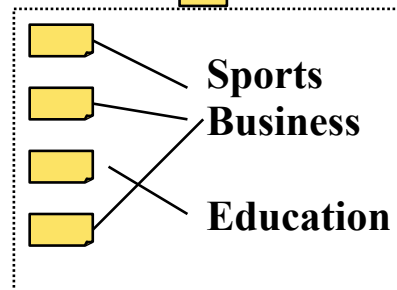
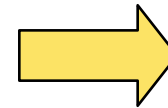
Price: **\$83.28**
Used & new from \$63.53

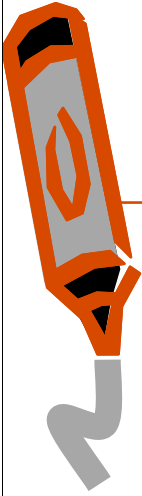
categorization

- Pre-given categories and labeled document examples (Categories may form hierarchy)
- Classify new documents
- A standard supervised learning problem



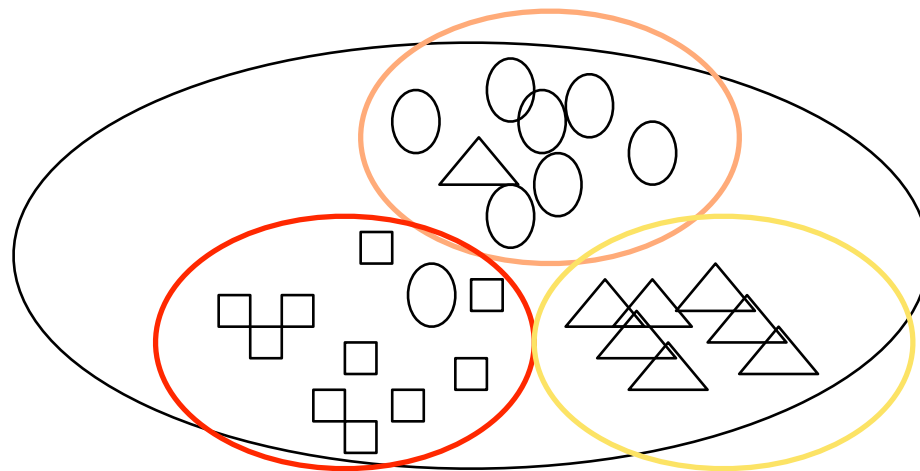
**Categorization
System**

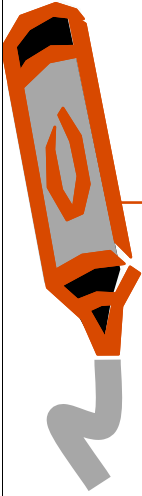




clustering

- Discover “natural structure”
- Group similar objects together
- Object can be document, term, passages

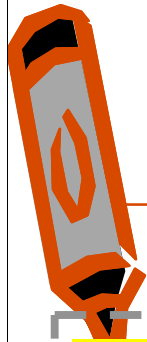




IR vs Databases

	Databases	IR
Data	Structured	Unstructured
Fields	Clear semantics (SSN, age)	No fields (other than text)
Queries	Defined (relational algebra, SQL)	Free text (“natural language”), Boolean
Recoverability	Critical (concurrency control, recovery, atomic operations)	Downplayed , though still an issue
Matching	Exact (results are <i>always</i> “correct”)	Imprecise (need to measure effectiveness)

related areas



Models

Applications

Applications
Web, Bioinformatics...

Machine Learning
Pattern Recognition
Data Mining

Library & Info
Science

Statistics
Optimization

**Information
Retrieval**

Databases

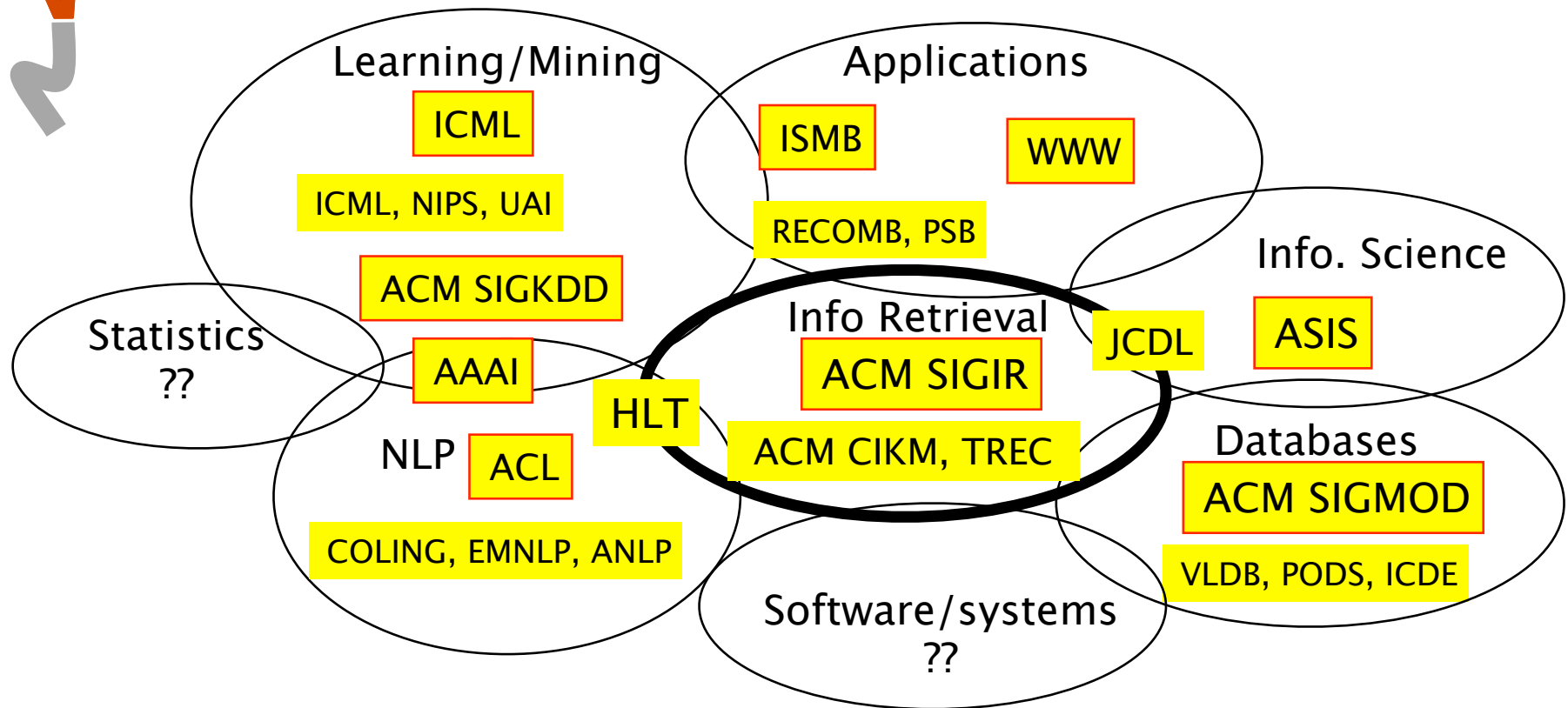
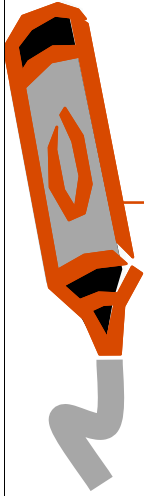
Natural
Language
Processing

Software engineering
Computer systems

Algorithms

Systems

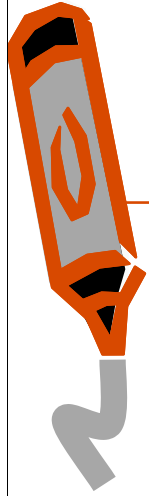
publications/societies





SIGIR 2003 topics

- Formal Models, Language Models, Fusion/Combination
- Text Representation and Indexing, XML and Metadata
- Performance, Compression, Scalability, Architectures, Mobile Applications
- Web IR, Intranet/Enterprise Search, Citation and Link Analysis, Digital Libraries, Distributed IR
- Cross-language Retrieval, Multilingual Retrieval, Machine Translation for IR
- Video and Image Access, Audio and Speech Retrieval, Music Retrieval
- Machine Learning for IR, Text Data Mining, Clustering, Text Categorization
- Topic Detection and Tracking, Content-Based Filtering, Collaborative Filtering, Agents
- Summarization, Question Answering, Natural Language Processing for IR, Information Extraction, Lexical Acquisition
- Interfaces, Visualization, Interactive IR, User Models, User Studies
- Specialized Applications of IR, including Genomic IR, IR in Software Engineering, and IR for Chemical Structures



overview

- what is information retrieval ?
- how does it work ?
- a simple retrieval model



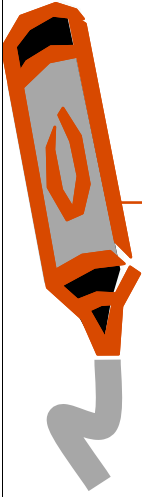
basic approaches

- boolean
- geometric : vector space model
- probabilistic : language models
- statistical : bayesian networks
- graph like : page rank



relevant items are similar

- Much of IR depends upon idea that **similar vocabulary → relevant to same queries**
- Usually look for documents matching query words
- “Similar” can be measured in many ways
 - String matching/comparison
 - Same vocabulary used
 - Probability that documents arise from same model
 - Same meaning of text

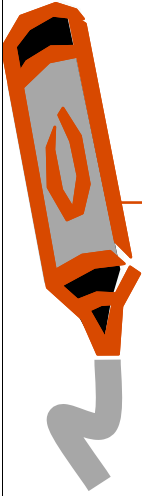


bag of words

- An effective and popular approach
 - Compares words without regard to order
 - Consider reordering words in a headline
-
- **Random:** beating takes points falling another Dow 355
 - **Alphabetical:** 355 another beating Dow falling points
 - **“Interesting”:** Dow points beating falling 355 another

 - **Actual: Dow takes another beating, falling 355 points**

what is this about ?



16 × said

12 × fat

8 × new

5 × food oil percent reduce taste Tuesday

4 × amount change health Henstenburg make obesity

3 × acids consumer fatty polyunsaturated US

2 × amounts artery Beemer cholesterol clogging director

down eat estimates expert fast formula impact initiative

moderate plans restaurant saturated trans win

1 × ...

added addition adults advocate affect afternoon age

Americans Asia battling beef bet brand Britt Brook Browns

calorie center chain chemically ... crispy customers cut ...

vegetable weapon weeks Wendys Wootan worldwide years

York

14 × McDonalds

11 × fries

6 × company french nutrition

the text



McDonald's slims down spuds

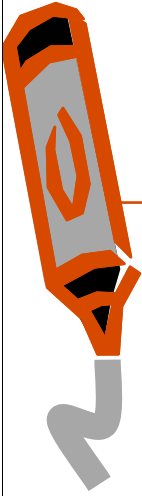
Fast-food chain to reduce certain types of fat in its french fries with new cooking oil.

NEW YORK (CNN/Money) - McDonald's Corp. is cutting the amount of "bad" fat in its french fries nearly in half, the fast-food chain said Tuesday as it moves to make all its fried menu items healthier.

But does that mean the popular shoestring fries won't taste the same? The company says no. "It's a win-win for our customers because they are getting the same great french-fry taste along with an even healthier nutrition profile," said Mike Roberts, president of McDonald's USA. But others are not so sure. McDonald's will not specifically discuss the kind of oil it plans to use, but at least one nutrition expert says playing with the formula could mean a different taste.

Shares of Oak Brook, Ill.-based McDonald's (MCD: down \$0.54 to \$23.22, Research, Estimates) were lower Tuesday afternoon. It was unclear Tuesday whether competitors Burger King and Wendy's International (WEN: down \$0.80 to \$34.91, Research, Estimates) would follow suit.

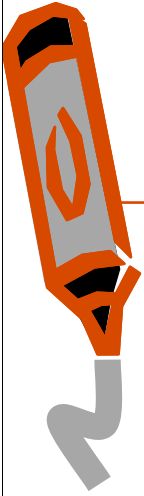
Neither company could immediately be reached for comment.



text representation

- Text representation
 - what makes a “good” representation?
 - how is a representation generated from text?
 - what are retrievable objects and how are they organized?
- Representing information needs
 - what is an appropriate query language?
 - how can interactive query formulation and refinement be supported?
- Comparing representations
 - what is a “good” model of retrieval?
 - how is uncertainty represented?

hypertext



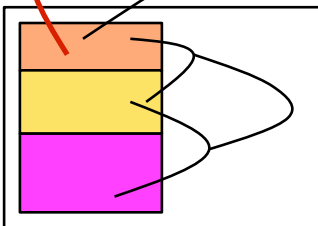
Concept map

A general topic

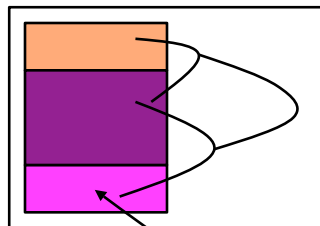
Subtopic 1

Subtopic i

Subtopic M

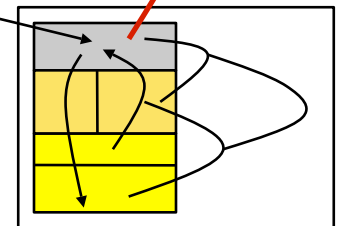
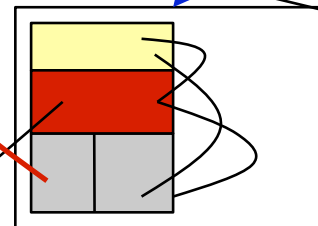


Doc 1



Doc 2

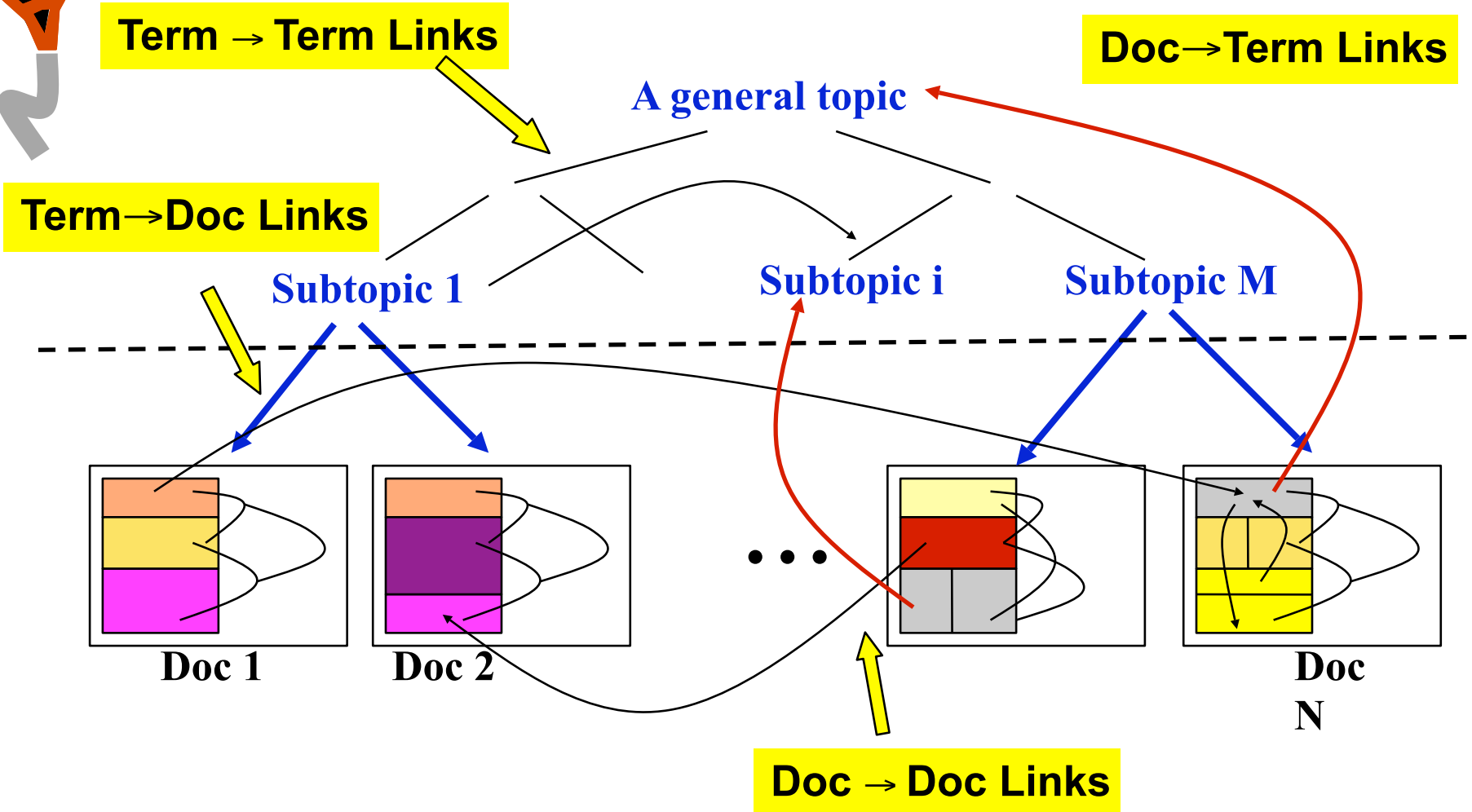
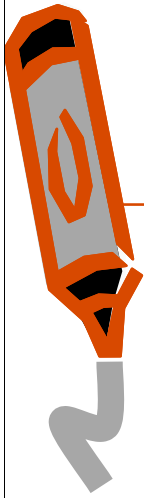
...



Doc
N

Hypertext

hypertext

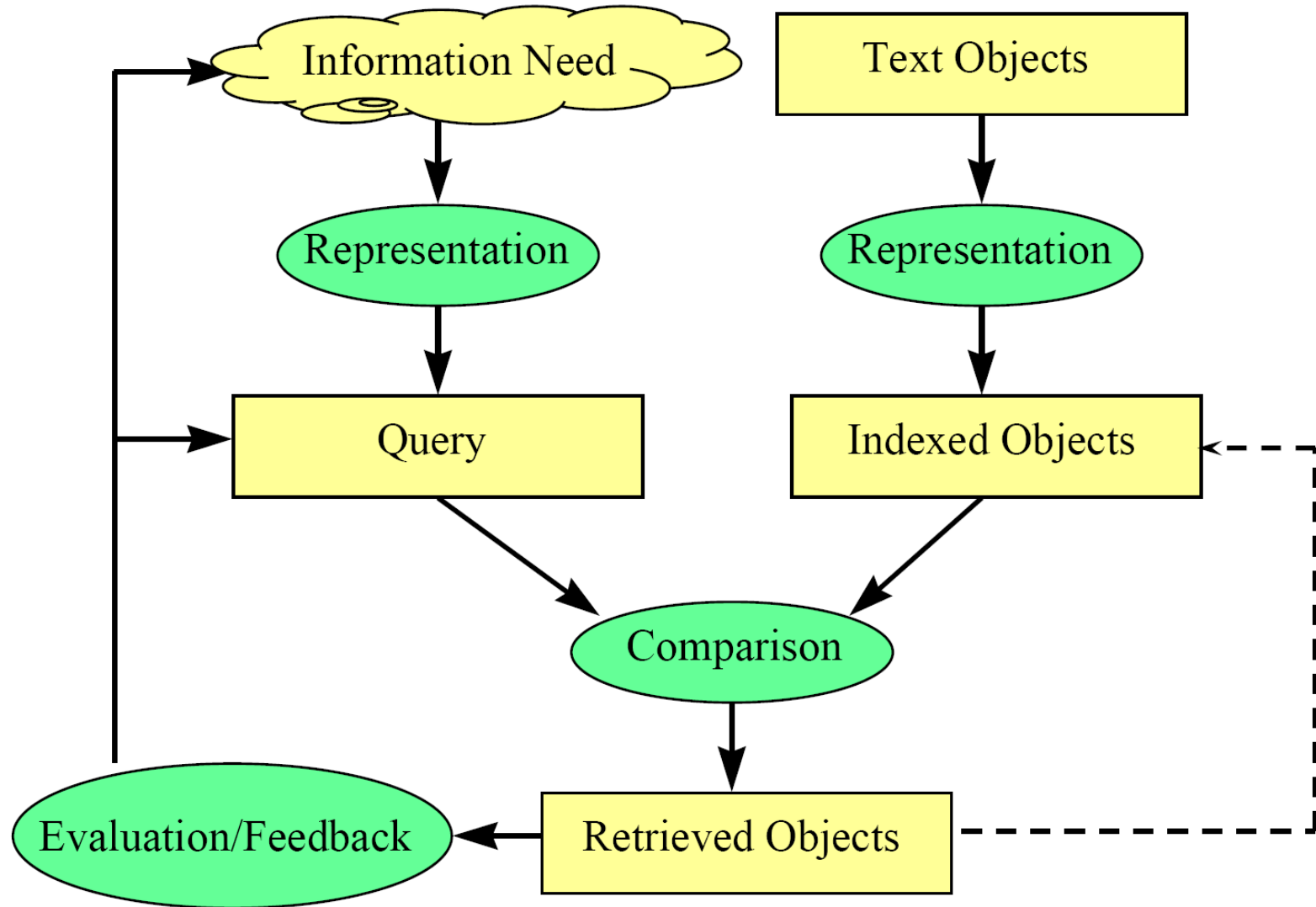


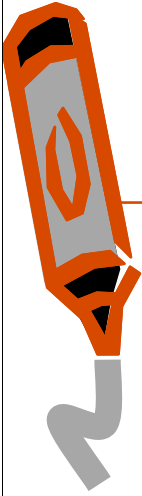


overview

- what is information retrieval ?
- how does it work ?
- a simple retrieval model

retrieval process





statistical language model

$D = \left\{ \begin{array}{l} \text{One fish, two fish, red fish, blue fish.} \\ \text{Black fish, blue fish, old fish, new fish.} \end{array} \right.$

$$\text{len}(D) = 16$$

$$P(\text{fish}|D) = 8/16 = 0.5$$

$$P(\text{blue}|D) = 2/16 = 0.125$$

$$P(\text{one}|D) = 1/16 = 0.0625$$

...

$$P(\text{eggs}|D) = 0/16 = 0$$

...


} A "topic"



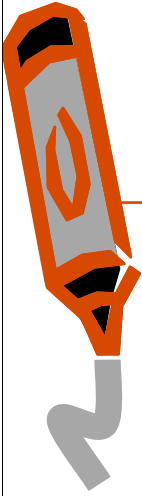
statistical language model

- Document came from a topic
- Did query come from *this* document's topic?

- For each document, find probability its topic could have generated the query

$$\begin{aligned} P(Q|T_D) &\approx P(Q|D) \\ &= P(q_1, \dots, q_t|D) \\ &= \prod_{i=1}^t P(q_i|D) \end{aligned}$$


Independence assumption
(Naïve Bayes)



statistical language model

$D_1 = \left\{ \begin{array}{l} \text{This one, I think, is called a Yink.} \\ \text{He likes to wink, he likes to drink.} \end{array} \right.$

$D_2 = \left\{ \begin{array}{l} \text{He likes to drink, and drink, and drink.} \\ \text{The thing he likes to drink is ink.} \end{array} \right.$

$D_3 = \left\{ \begin{array}{l} \text{The ink he likes to drink is pink.} \\ \text{He links to wink and drink pink ink.} \end{array} \right.$

Query “drink”

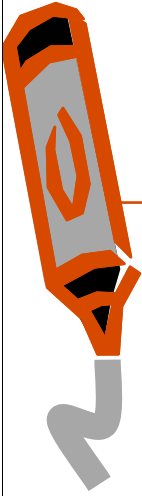
- $P(\text{drink}|D_1) = 1/16$
- $P(\text{drink}|D_2) = 4/16$
- $P(\text{drink}|D_3) = 2/16$

Query “pink ink”

- $P(Q|D_1) = 0 \cdot 0 = 0$
- $P(Q|D_2) = 0 \cdot 1/16 = 0$
- $P(Q|D_3) = 2/16 \cdot 2/16 = 0.016$

Query “wink drink”

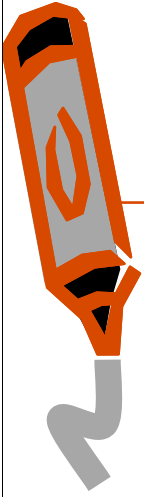
- $P(Q|D_1) = 0.004$
- $P(Q|D_2) = 0$
- $P(Q|D_3) = 1/16 \cdot 2/16 = 0.008$



conclusion

- Information Retrieval?
 - Indexing, retrieving, and organizing text by probabilistic or statistical techniques that reflect semantics without actually understanding
 - Search engines
- Core idea
 - Bag of words captures much of the “meaning”
 - Objects that use vocabulary the same way are related
- Statistical language model
 - Documents used to estimate a topic model
 - Query reflects a topic, too
 - Documents of topics that are likely to produce the query are most likely to be relevant

in this course



- <http://ccs.neu.edu/~jaa/ISU535.06X2/schedulen.html>