# Text Retrieval and Web Search
# IIR 1: Boolean Retrieval

Mihai Surdeanu

(Based on slides by Hinrich Schütze at `informationretrieval.org`)

Spring 2017

## Take-away

- Why you should take this course
- Admin issues
- Boolean Retrieval: Design and data structures of a simple information retrieval system
- What topics will be covered in this class?

# Why you should take this course (take 1)

# Why you should take this course (take 2)



"Making the world's problem solvers 10% more efficient"

# Why you should take this course (take 2)

- Anurag Acharya is the key inventor of Google Scholar.
- `https://medium.com/backchannel/`
  `the-gentleman-who-made-scholar-d71289d9a82d`

## Outline

1 **Administration**

2 Introduction

3 Inverted index

4 Processing Boolean queries

5 Query optimization

6 Course overview

## Instructor information

Instructor: Mihai Surdeanu
Email: `msurdeanu@email.arizona.edu`
Office: Gould-Simpson 746
Office Hours: Mon/Wed 11 – noon

TA: Enrique Noriega
Email: `enoriega@email.arizona.edu`
Office: Gould-Simpson 931
Office Hours: Tue/Thu 11 – 12:30

## Website and syllabus

- Website:
  http://surdeanu.info/mihai/teaching/
  csc483583-spring17/
- Syllabus:
  http://surdeanu.info/mihai/teaching/
  csc483583-spring17/IR-syllabus.pdf
- Piazza:
  https:
  //piazza.com/arizona/spring2017/csc483583/home
- See website and syllabus for: code of conduct, classroom electronics, DRC accommodations, non-discrimination policy, code of academic integrity, etc.
- But most materials will actually be in D2L

## Prerequisites

- Know how to program and have a decent understanding of data structures such as hash maps and trees: CSC 345
- Ideally, Math 129 (Calc 2). However, we will cover the necessary math in class.

# Prerequisites: does this look scary?

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

```
 1   answer ← ⟨ ⟩
 2   while p₁ ≠ NIL and p₂ ≠ NIL
 3   do if docID(p₁) = docID(p₂)
 4         then ADD(answer, docID(p₁))
 5               p₁ ← next(p₁)
 6               p₂ ← next(p₂)
 7         else if docID(p₁) < docID(p₂)
 8               then if hasSkip(p₁) and (docID(skip(p₁)) ≤ docID(p₂))
 9                     then while hasSkip(p₁) and (docID(skip(p₁)) ≤ docID(p₂))
10                           do p₁ ← skip(p₁)
11                     else  p₁ ← next(p₁)
12               else if hasSkip(p₂) and (docID(skip(p₂)) ≤ docID(p₁))
13                     then while hasSkip(p₂) and (docID(skip(p₂)) ≤ docID(p₁))
14                           do p₂ ← skip(p₂)
15                     else  p₂ ← next(p₂)
16   return answer
```

# Prerequisites: does this look scary?

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

$$\cos(\vec{q}, \vec{d}) = \text{SIM}(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{|\vec{q}||\vec{d}|} = \frac{\sum_{i=1}^{|V|} q_i d_i}{\sqrt{\sum_{i=1}^{|V|} q_i^2} \sqrt{\sum_{i=1}^{|V|} d_i^2}}$$

Dot product, matrix multiplication, Bayes rule

# Choosing a programming language
## My recommendations

- Scala
- Java
- Python
- C/C++

## Java

- Pros
    - Pretty fast
    - Probably the most common language for IR and NLP
    - Clean exception handling
    - Statically typed
    - Garbage collection
    - Several great IDEs

- Cons
    - Syntax too verbose
    - Inconsistent semantics due to enforced backwards compatibility
      (primitive types vs. objects, equality, etc.)

## Scala

- Pros
    - Pretty fast
    - "Hot" language for IR, NLP, ML, distributed computing, web development
    - Clean, transparent exception handling
    - Clean, minimalist syntax
    - Consistent semantics
    - Statically typed
    - Garbage collection
    - At least one great IDE (IntelliJ)
    - Fully compatible with Java (use all Java libraries)

- Cons
    - It has some "dark corners"
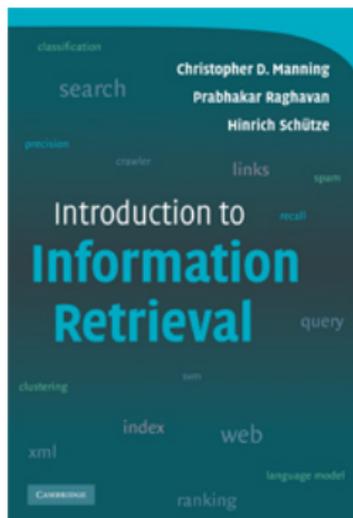    - Backwards compatibility not guaranteed

# Python

- Pros
  - Clean syntax
  - Popular: many NLP/ML libraries exist
  - Clean exception handling
  - Easy access to GPUs

- Cons
  - **Slow**
  - **Dynamically typed**
  - No great IDE

# C/C++

- Pros
  - As fast as it gets
- Cons
  - Too many to list

# Comparison



Ten tiny examples - How many times slower?

More benchmarks:

http://benchmarksgame.alioth.debian.org/u64/benchmark.php?test=all&lang=all&data=u64

## Textbook



- Introduction to Information Retrieval, by Manning et al.
- http://nlp.stanford.edu/IR-book/

# Grading

- Final grade = 4 assignments + 2 exams + 1 project + in-class participation
- First assignment has been posted. Due January 29! Let's take a quick look at the assignment.
- Undergraduate vs. graduate requirements

# Grading scheme

| Component | Weight |
|---|---|
| Written assignments | 300 pts |
| Midterm exam | 200 pts |
| Final exam | 275 pts |
| Programming project | 200 pts |
| In-class participation | 25 pts |
| Total | 1000 pts |

| Grade | Point Range |
|---|---|
| A | 900 − 1000 |
| B | 800 − 899 |
| C | 700 − 799 |
| D | 600 − 699 |
| E | 0 − 599 |

# Final Project: Option 1

# Final Project: Option 2

**Completed • $80,000 • 170 teams**

## The Allen AI Science Challenge

Wed 7 Oct 2015 – Sat 13 Feb 2016 (11 months ago)

**ALLEN INSTITUTE**
*for* **ARTIFICIAL INTELLIGENCE**

Competition Details  »  Get the Data  »  Make a submission

## Is your model smarter than an 8th grader?

**Dashboard**

Home
Data
Make a submission

Information
Description
Evaluation
Rules
Prizes
Timeline

Forum

Leaderboard
Public
Private

Private Leaderboard

The Allen Institute for Artificial Intelligence (AI2) is working to improve humanity through fundamental advances in artificial intelligence. One critical but challenging problem in AI is to demonstrate the ability to consistently understand and correctly answer general questions about the world.

# Final Project: Option 3

## Post-facto Fake News Challenge

Register a team     Our Github repositories     Join the Slack

## Description

Post-facto fake news refers to news items or claims that are already known to be false, either by work from organizations like Snopes and Factcheck or by the general public on social media.

For this challenge, we will only consider claims that are outright false or outright true. For example, "eating fruit prevents cancer" -- a truth assertion here is dodgy at best, but for headlines like "Hillary has a body double" we can be confident about truth assertion.

We will select headlines for the competition where we can be confident in asserting its veracity.

# Late work + attendance policy

- Late work is not accepted, except in case of documented emergency approved by the instructor

- Attendance is required
- Students who miss class due to illness or emergency are required to bring documentation

## Dates

| What | When |
|---|---|
| HW 1 | January 29 |
| HW 2 | February 19 |
| Midterm | March 1 |
| HW 3 | March 12 |
| Spring recess | March 13 - 17 |
| HW 4 | April 9 |
| Project | May 7 |
| Final | After last week of class |

Written assignments will be due approximately every three weeks, as announced by the instructor. All assignments are due in the D2L dropbox by 11:59 P.M. on the indicated day.

## Cooperation & Cheating

- Students may not discuss individual homework with anybody other than the instructors and teaching assistants.
- Students may not share individual homework solutions with anybody.
- Students may post questions to Piazza, but should refrain from posting solutions or partial solutions.
- Students may not share test cases with anybody.
- Students may share class notes with anybody.
- Students may not seek individual homework help from anybody other than the instructors, teaching assistants, or departmental tutors.
- If permitted, the use of open source or third party materials in student submissions must be clearly identified and credited. Assignment and project submissions must be substantially the work of the student who submits the work.

## Cooperation & Cheating

- Students who violate the Code should expect a penalty that is **greater than the value of the work in question up to and including failing the course**.
- A record of the incident **will** be sent to the Dean of Students office. If you have been involved in other Code violations, the Dean of Students may impose additional sanctions.

# Outline

## Definition of *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).

# Definition of *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).

# Definition of *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).

# Definition of *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).

# Definition of *information retrieval*

Information retrieval (IR) is finding material (usually documents) of
an unstructured nature (usually text) that satisfies an information
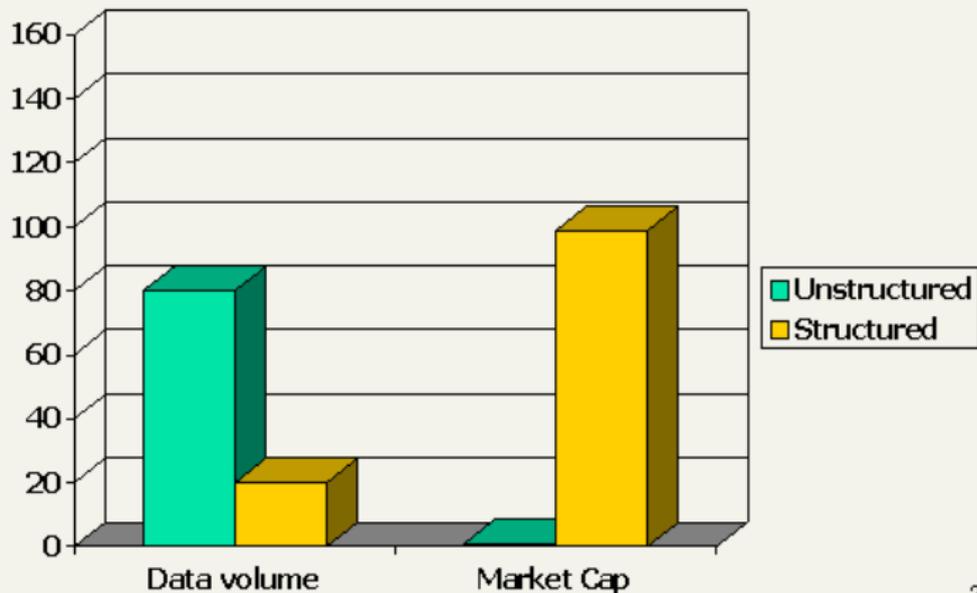need from within large collections (usually stored on computers).

# Definition of *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).

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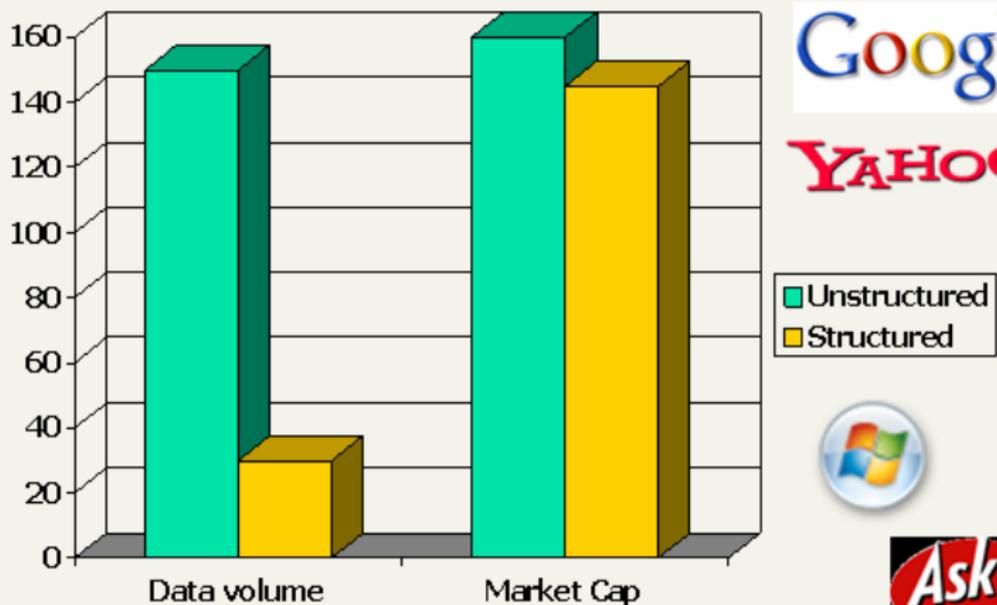
# Unstructured (text) vs. structured (database) data in 1996

# Unstructured (text) vs. structured (database) data in 2006

# Companies using IR

- Google, Yahoo, Microsoft: search web, email, choose ads
- Facebook: search friends' posts, choose wall
- Twitter: search tweets
- HP Autonomy: enterprise search
- Pandora: music (!) search

# Why you should take this course

- Take 3: IR is important

- Take 4:
    - This is a gateway course for data science and machine learning
    - Real-world applications of programming + data structures + probability theory

# Boolean retrieval

- The Boolean model is arguably the simplest model to base an information retrieval system on.
- Queries are Boolean expressions, e.g., CAESAR AND BRUTUS
- The search engine returns all documents that satisfy the Boolean expression.

## Boolean retrieval

- The Boolean model is arguably the simplest model to base an information retrieval system on.
- Queries are Boolean expressions, e.g., CAESAR AND BRUTUS
- The search engine returns all documents that satisfy the Boolean expression.

Does Google use the Boolean model?

## Does Google use the Boolean model?

- On Google, the default interpretation of a query [$w_1$ $w_2$ ... $w_n$] is $w_1$ AND $w_2$ AND ... AND $w_n$
- Cases where you get hits that do not contain one of the $w_i$:
  - anchor text
  - page contains variant of $w_i$ (morphology, spelling correction, synonym)
  - long queries ($n$ large)
  - boolean expression generates very few hits
- Simple Boolean vs. Ranking of result set
  - Simple Boolean retrieval returns matching documents in no particular order.
  - Google (and most well designed Boolean engines) rank the result set – they rank good hits (according to some estimator of relevance) higher than bad hits.

# Outline

# Unstructured data in 1650: Shakespeare

# Unstructured data in 1650

- 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.
- Why is grep not the solution?

## Unstructured data in 1650

- 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.
- Why is grep not the solution?
  - Slow (for large collections)
  - grep is line-oriented, IR is document-oriented
  - "NOT CALPURNIA" is non-trivial
  - Other operations (e.g., find the word ROMANS near COUNTRYMAN) not feasible

## Term-document incidence matrix

|  | Anthony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth | ... |
|---|---|---|---|---|---|---|---|
| ANTHONY | 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 | |

...

Entry is 1 if term occurs. Example: CALPURNIA occurs in *Julius Caesar*.
Entry is 0 if term doesn't occur. Example: CALPURNIA doesn't occur in *The tempest*.

# Term-document incidence matrix

|  | Anthony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth | . . . |
|---|---|---|---|---|---|---|---|
| ANTHONY | 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 | |

. . .

Entry is 1 if term occurs. Example: CALPURNIA occurs in *Julius Caesar*.

Entry is 0 if term doesn't occur. Example: CALPURNIA doesn't occur in *The tempest*.

# Term-document incidence matrix

| | Anthony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth | . . . |
|---|---|---|---|---|---|---|---|
| Anthony | 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 | |

. . .

Entry is 1 if term occurs. Example: Calpurnia occurs in *Julius Caesar*.
Entry is 0 if term doesn't occur. Example: Calpurnia doesn't occur in *The tempest*.

## Incidence vectors

- So we have a 0/1 vector for each term.
- To answer the query BRUTUS AND CAESAR AND NOT CALPURNIA:

## Incidence vectors

- So we have a 0/1 vector for each term.
- To answer the query BRUTUS AND CAESAR AND NOT CALPURNIA:
  - Take the vectors for BRUTUS, CAESAR, and CALPURNIA
  - Complement the vector of CALPURNIA
  - Do a (bitwise) AND on the three vectors
  - 110100 AND 110111 AND 101111 = 100100

# 0/1 vector for BRUTUS

| | Anthony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth | ... |
|---|---|---|---|---|---|---|---|
| ANTHONY | 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 | |
| ... | | | | | | | |
| result: | 1 | 0 | 0 | 1 | 0 | 0 | |

## Answers to query

*Anthony 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.

# Bigger collections

- Consider $N = 10^6$ documents, each with about 1000 tokens
- $\Rightarrow$ total of $10^9$ tokens
- On average 6 bytes per token, including spaces and punctuation $\Rightarrow$ size of document collection is about $6 \cdot 10^9 = 6$ GB
- Assume there are $M = 500{,}000$ distinct terms in the collection
- (Notice that we are making a term/token distinction.)

## Can't build the incidence matrix

- $M = 500{,}000 \times 10^6 =$ half a trillion 0s and 1s.
- But the matrix has no more than one billion 1s.
  - Matrix is extremely sparse.
- What is a better representations?
  - We only record the 1s.

## Inverted Index

For each term $t$, we store a list of all documents that contain $t$.



dictionary                  postings

# Inverted Index

For each term $t$, we store a list of all documents that contain $t$.

| BRUTUS | $\longrightarrow$ | 1 | 2 | 4 | 11 | 31 | 45 | 173 | 174 |
|---|---|---|---|---|---|---|---|---|---|

| CAESAR | $\longrightarrow$ | 1 | 2 | 4 | 5 | 6 | 16 | 57 | 132 | ... |
|---|---|---|---|---|---|---|---|---|---|---|

| CALPURNIA | $\longrightarrow$ | 2 | 31 | 54 | 101 |
|---|---|---|---|---|---|

⋮

**dictionary**　　　　　　　　　　**postings**

# Inverted Index

For each term $t$, we store a list of all documents that contain $t$.

| Brutus | $\longrightarrow$ | 1 | 2 | 4 | 11 | 31 | 45 | 173 | 174 |
|---|---|---|---|---|---|---|---|---|---|

| Caesar | $\longrightarrow$ | 1 | 2 | 4 | 5 | 6 | 16 | 57 | 132 | ... |
|---|---|---|---|---|---|---|---|---|---|---|

| Calpurnia | $\longrightarrow$ | 2 | 31 | 54 | 101 |
|---|---|---|---|---|---|

$\vdots$

$\underbrace{\phantom{xxxxxxxxx}}_{\textbf{dictionary}}$     $\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}}_{\textbf{postings}}$

# Inverted index construction

1. Collect the documents to be indexed:

   | Friends, Romans, countrymen. | | So let it be with Caesar | . . .

2. Tokenize the text, turning each document into a list of tokens:

   | Friends | | Romans | | countrymen | | So | . . .

3. Do linguistic preprocessing, producing a list of normalized tokens, which are the indexing terms: | friend | | roman |

   | countryman | | so | . . .

4. Index the documents that each term occurs in by creating an inverted index, consisting of a dictionary and postings.

# Tokenization and preprocessing

**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:

$\Longrightarrow$

**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

# Generate postings

| term | docID |
|------|-------|
| 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 |

**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

$\Longrightarrow$

# Sort postings

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

# Create postings lists, determine document frequency

| 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 |

$\Longrightarrow$

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

# Split the result into dictionary and postings file

| Brutus | $\longrightarrow$ | 1 | 2 | 4 | 11 | 31 | 45 | 173 | 174 |
|---|---|---|---|---|---|---|---|---|---|

| Caesar | $\longrightarrow$ | 1 | 2 | 4 | 5 | 6 | 16 | 57 | 132 | ... |
|---|---|---|---|---|---|---|---|---|---|---|

| Calpurnia | $\longrightarrow$ | 2 | 31 | 54 | 101 |
|---|---|---|---|---|---|

⋮

**dictionary**                          **postings file**

## Later in this course

- Index construction: how can we create inverted indexes for large collections?
- How much space do we need for dictionary and index?
- Index compression: how can we efficiently store and process indexes for large collections?
- Ranked retrieval: what does the inverted index look like when we want the "best" answer?

# Outline

# Simple conjunctive query (two terms)

- Consider the query: BRUTUS AND CALPURNIA
- To find all matching documents using inverted index:
  1. Locate BRUTUS in the dictionary
  2. Retrieve its postings list from the postings file
  3. Locate CALPURNIA in the dictionary
  4. Retrieve its postings list from the postings file
  5. Intersect the two postings lists
  6. Return intersection to user

# Intersecting two postings lists

BRUTUS $\longrightarrow$ $\boxed{1} \to \boxed{2} \to \boxed{4} \to \boxed{11} \to \boxed{31} \to \boxed{45} \to \boxed{173} \to \boxed{174}$

CALPURNIA $\longrightarrow$ $\boxed{2} \to \boxed{31} \to \boxed{54} \to \boxed{101}$

Intersection $\implies$

# Intersecting two postings lists

BRUTUS $\longrightarrow$ $\boxed{1} \to \boxed{2} \to \boxed{4} \to \boxed{11} \to \boxed{31} \to \boxed{45} \to \boxed{173} \to \boxed{174}$

CALPURNIA $\longrightarrow$ $\boxed{2} \to \boxed{31} \to \boxed{54} \to \boxed{101}$

Intersection $\implies$

# Intersecting two postings lists

BRUTUS $\longrightarrow$ $\boxed{1} \to \boxed{2} \to \boxed{4} \to \boxed{11} \to \boxed{31} \to \boxed{45} \to \boxed{173} \to \boxed{174}$

CALPURNIA $\longrightarrow$ $\boxed{2} \to \boxed{31} \to \boxed{54} \to \boxed{101}$

Intersection $\Longrightarrow$

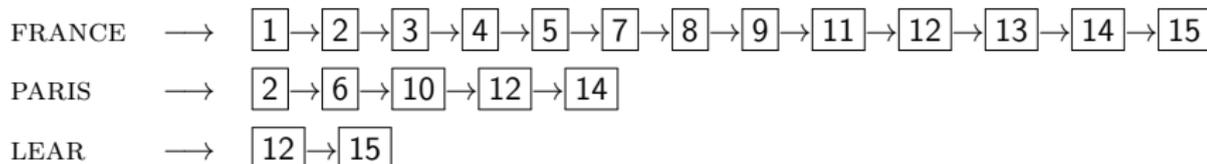# Intersecting two postings lists

BRUTUS $\longrightarrow$ $\boxed{1} \rightarrow \boxed{2} \rightarrow \boxed{4} \rightarrow \boxed{11} \rightarrow \boxed{31} \rightarrow \boxed{45} \rightarrow \boxed{173} \rightarrow \boxed{174}$

CALPURNIA $\longrightarrow$ $\boxed{2} \rightarrow \boxed{31} \rightarrow \boxed{54} \rightarrow \boxed{101}$

Intersection $\implies$ $\boxed{2}$

# Intersecting two postings lists

BRUTUS $\longrightarrow$ $\boxed{1} \rightarrow \boxed{2} \rightarrow \boxed{4} \rightarrow \boxed{11} \rightarrow \boxed{31} \rightarrow \boxed{45} \rightarrow \boxed{173} \rightarrow \boxed{174}$

CALPURNIA $\longrightarrow$ $\boxed{2} \rightarrow \boxed{31} \rightarrow \boxed{54} \rightarrow \boxed{101}$

Intersection $\implies$ $\boxed{2}$

# Intersecting two postings lists

Brutus         $\longrightarrow$    $\boxed{1} \to \boxed{2} \to \boxed{4} \to \boxed{11} \to \boxed{31} \to \boxed{45} \to \boxed{173} \to \boxed{174}$

Calpurnia   $\longrightarrow$    $\boxed{2} \to \boxed{31} \to \boxed{54} \to \boxed{101}$

Intersection  $\implies$    $\boxed{2}$

# Intersecting two postings lists

BRUTUS $\longrightarrow$ $\boxed{1} \rightarrow \boxed{2} \rightarrow \boxed{4} \rightarrow \boxed{11} \rightarrow \boxed{31} \rightarrow \boxed{45} \rightarrow \boxed{173} \rightarrow \boxed{174}$

CALPURNIA $\longrightarrow$ $\boxed{2} \rightarrow \boxed{31} \rightarrow \boxed{54} \rightarrow \boxed{101}$

Intersection $\implies$ $\boxed{2}$

# Intersecting two postings lists

| | | |
|---|---|---|
| BRUTUS | $\longrightarrow$ | $1 \to 2 \to 4 \to 11 \to 31 \to 45 \to 173 \to 174$ |
| CALPURNIA | $\longrightarrow$ | $2 \to 31 \to 54 \to 101$ |
| Intersection | $\implies$ | $2 \to 31$ |

# Intersecting two postings lists

BRUTUS $\longrightarrow$ $\boxed{1} \to \boxed{2} \to \boxed{4} \to \boxed{11} \to \boxed{31} \to \boxed{45} \to \boxed{173} \to \boxed{174}$

CALPURNIA $\longrightarrow$ $\boxed{2} \to \boxed{31} \to \boxed{54} \to \boxed{101}$

Intersection $\implies$ $\boxed{2} \to \boxed{31}$

# Intersecting two postings lists

BRUTUS $\longrightarrow$ $\boxed{1} \to \boxed{2} \to \boxed{4} \to \boxed{11} \to \boxed{31} \to \boxed{45} \to \boxed{173} \to \boxed{174}$

CALPURNIA $\longrightarrow$ $\boxed{2} \to \boxed{31} \to \boxed{54} \to \boxed{101}$

Intersection $\implies$ $\boxed{2} \to \boxed{31}$

# Intersecting two postings lists

BRUTUS $\longrightarrow$ $\boxed{1} \rightarrow \boxed{2} \rightarrow \boxed{4} \rightarrow \boxed{11} \rightarrow \boxed{31} \rightarrow \boxed{45} \rightarrow \boxed{173} \rightarrow \boxed{174}$

CALPURNIA $\longrightarrow$ $\boxed{2} \rightarrow \boxed{31} \rightarrow \boxed{54} \rightarrow \boxed{101}$

Intersection $\implies$ $\boxed{2} \rightarrow \boxed{31}$

## Intersecting two postings lists

BRUTUS $\longrightarrow$ $\boxed{1} \rightarrow \boxed{2} \rightarrow \boxed{4} \rightarrow \boxed{11} \rightarrow \boxed{31} \rightarrow \boxed{45} \rightarrow \boxed{173} \rightarrow \boxed{174}$

CALPURNIA $\longrightarrow$ $\boxed{2} \rightarrow \boxed{31} \rightarrow \boxed{54} \rightarrow \boxed{101}$

Intersection $\implies$ $\boxed{2} \rightarrow \boxed{31}$

# Intersecting two postings lists

BRUTUS $\longrightarrow$ $\boxed{1} \rightarrow \boxed{2} \rightarrow \boxed{4} \rightarrow \boxed{11} \rightarrow \boxed{31} \rightarrow \boxed{45} \rightarrow \boxed{173} \rightarrow \boxed{174}$

CALPURNIA $\longrightarrow$ $\boxed{2} \rightarrow \boxed{31} \rightarrow \boxed{54} \rightarrow \boxed{101}$

Intersection $\implies$ $\boxed{2} \rightarrow \boxed{31}$

- This is linear in the length of the postings lists.

# Intersecting two postings lists

BRUTUS        $\longrightarrow$   $\boxed{1}\to\boxed{2}\to\boxed{4}\to\boxed{11}\to\boxed{31}\to\boxed{45}\to\boxed{173}\to\boxed{174}$

CALPURNIA   $\longrightarrow$   $\boxed{2}\to\boxed{31}\to\boxed{54}\to\boxed{101}$

Intersection  $\implies$   $\boxed{2}\to\boxed{31}$

- This is linear in the length of the postings lists.
- Note: This only works if postings lists are sorted.

# Intersecting two postings lists

Intersect($p_1, p_2$)
  1   *answer* ← ⟨ ⟩
  2   **while** $p_1 \neq$ NIL and $p_2 \neq$ NIL
  3   **do if** $docID(p_1) = docID(p_2)$
  4         **then** 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*

# Query processing: Exercise

FRANCE $\longrightarrow$ $\boxed{1} \rightarrow \boxed{2} \rightarrow \boxed{3} \rightarrow \boxed{4} \rightarrow \boxed{5} \rightarrow \boxed{7} \rightarrow \boxed{8} \rightarrow \boxed{9} \rightarrow \boxed{11} \rightarrow \boxed{12} \rightarrow \boxed{13} \rightarrow \boxed{14} \rightarrow \boxed{15}$

PARIS $\longrightarrow$ $\boxed{2} \rightarrow \boxed{6} \rightarrow \boxed{10} \rightarrow \boxed{12} \rightarrow \boxed{14}$

LEAR $\longrightarrow$ $\boxed{12} \rightarrow \boxed{15}$

Compute hit list for ((paris AND NOT france) OR lear)

## Boolean queries

- The Boolean retrieval model can answer any query that is a Boolean expression.
  - Boolean queries are queries that use AND, OR and NOT to join query terms.
  - Views each document as a set of terms.
  - Is precise: Document matches condition or not.
- Primary commercial retrieval tool for 3 decades
- Many professional searchers (e.g., lawyers) still like Boolean queries.
  - You know exactly what you are getting.
- Many search systems you use are also Boolean: spotlight, email, intranet etc.

# Commercially successful Boolean retrieval: Westlaw

- Largest commercial legal search service in terms of the number of paying subscribers
- Over half a million subscribers performing millions of searches a day over tens of terabytes of text data
- The service was started in 1975.
- In 2005, Boolean search (called "Terms and Connectors" by Westlaw) was still the default, and used by a large percentage of users . . .
- . . . although ranked retrieval has been available since 1992. Why do you think this is the case?

# Westlaw: Example queries

*Information need:* Information on the legal theories involved in preventing the disclosure of trade secrets by employees formerly employed by a competing company

*Query:* "trade secret" /s disclos! /s prevent /s employe!

# Westlaw: Example queries

*Information need:* Requirements for disabled people to be able to access a workplace

*Query:* disab! /p access! /s work-site work-place (employment /3 place)

# Westlaw: Example queries

*Information need:* Cases about a host's responsibility for drunk guests

*Query:* host! /p (responsib! liab!) /p (intoxicat! drunk!) /p guest

# Westlaw: Comments

- Proximity operators: /3 = within 3 words, /s = within a sentence, /p = within a paragraph
- Space is disjunction, not conjunction! (This was the default in search pre-Google.) Why do you think this was the case?
- Long, precise queries: incrementally developed, not like web search
- Why professional searchers often like Boolean search: precision, transparency, control
- When are Boolean queries the best way of searching? Depends on: information need, searcher, document collection, . . .

# Outline

## Query optimization

- Consider a query that is an AND of $n$ terms, $n > 2$
- For each of the terms, get its postings list, then AND them together
- Example query: BRUTUS AND CALPURNIA AND CAESAR
- What is the best order for processing this query?

## Query optimization

- Example query: BRUTUS AND CALPURNIA AND CAESAR

# Query optimization

- Example query: BRUTUS AND CALPURNIA AND CAESAR
- Simple and effective optimization: Process in order of increasing frequency. Why?

# Query optimization

- Example query: BRUTUS AND CALPURNIA AND CAESAR
- Simple and effective optimization: Process in order of increasing frequency. Why?
- Start with the shortest postings list, then keep cutting further

# Query optimization

- Example query: BRUTUS AND CALPURNIA AND CAESAR
- Simple and effective optimization: Process in order of increasing frequency. Why?
- Start with the shortest postings list, then keep cutting further

BRUTUS      $\longrightarrow$      $\boxed{1} \to \boxed{2} \to \boxed{4} \to \boxed{11} \to \boxed{31} \to \boxed{45} \to \boxed{173} \to \boxed{174}$

CALPURNIA   $\longrightarrow$      $\boxed{2} \to \boxed{31} \to \boxed{54} \to \boxed{101}$

CAESAR      $\longrightarrow$      $\boxed{5} \to \boxed{31}$

## Query optimization

- Example query: BRUTUS AND CALPURNIA AND CAESAR
- Simple and effective optimization: Process in order of increasing frequency. Why?
- Start with the shortest postings list, then keep cutting further
- In this example, first CAESAR, then CALPURNIA, then BRUTUS

BRUTUS     $\longrightarrow$    $\boxed{1} \to \boxed{2} \to \boxed{4} \to \boxed{11} \to \boxed{31} \to \boxed{45} \to \boxed{173} \to \boxed{174}$

CALPURNIA  $\longrightarrow$    $\boxed{2} \to \boxed{31} \to \boxed{54} \to \boxed{101}$

CAESAR     $\longrightarrow$    $\boxed{5} \to \boxed{31}$

# Optimized intersection algorithm for conjunctive queries

$\text{Intersect}(\langle t_1, \ldots, t_n \rangle)$
1   $terms \leftarrow \text{SortByIncreasingFrequency}(\langle t_1, \ldots, t_n \rangle)$
2   $result \leftarrow postings(first(terms))$
3   $terms \leftarrow rest(terms)$
4   **while** $terms \neq \text{NIL}$ and $result \neq \text{NIL}$
5   **do** $result \leftarrow \text{Intersect}(result, postings(first(terms)))$
6       $terms \leftarrow rest(terms)$
7   **return** $result$

# More general optimization

- Example query: (MADDING OR CROWD) AND (IGNOBLE OR STRIFE)
- Get frequencies for all terms
- Estimate the size of each OR by the sum of its frequencies (conservative)
- Process in increasing order of OR sizes

## Outline

## Course overview

- We are done with Chapter 1 of IIR (IIR 01).
- Plan for the rest of the semester: 18–20 of the 21 chapters of IIR
- In what follows: teasers for most chapters – to give you a sense of what will be covered.

# IIR 02: The term vocabulary and postings lists

- Phrase queries: "STANFORD UNIVERSITY"
- Proximity queries: GATES NEAR MICROSOFT
- We need an index that captures position information for phrase queries and proximity queries.

# IIR 03: Dictionaries and tolerant retrieval

Wildcard queries, spelling correction

| BO | → | aboard | → | about | → | boardroom | → | border |
|---|---|---|---|---|---|---|---|---|

| OR | → | border | → | lord | → | morbid | → | sordid |
|---|---|---|---|---|---|---|---|---|

| RD | → | aboard | → | ardent | → | boardroom | → | border |
|---|---|---|---|---|---|---|---|---|

# IIR 04: Index construction

# IIR 05: Index compression

Zipf's law:

# IIR 06: The vector space model

- Ranking search results
  - Boolean queries only give inclusion or exclusion of documents.
  - For ranked retrieval, we measure the proximity between the query and each document.
  - One formalism for doing this: the vector space model
- Key challenge in ranked retrieval: evidence accumulation for a term in a document
  - 1 vs. 0 occurence of a query term in the document
  - 3 vs. 2 occurences of a query term in the document
  - Usually: more is better
  - But by how much?
  - Need a scoring function that translates frequency into score or weight

# IIR 07: Scoring in a complete search system

# IIR 08: Evaluation and dynamic summaries

# IIR 09: Relevance feedback & query expansion

Keeping the human in the loop:

# IIR 11: Probabilistic information retrieval

$P(R|d, q)$ is modeled using term incidence vectors as $P(R|\vec{x}, \vec{q})$

$$P(R = 1|\vec{x}, \vec{q}) = \frac{P(\vec{x}|R = 1, \vec{q})P(R = 1|\vec{q})}{P(\vec{x}|\vec{q})}$$

$$P(R = 0|\vec{x}, \vec{q}) = \frac{P(\vec{x}|R = 0, \vec{q})P(R = 0|\vec{q})}{P(\vec{x}|\vec{q})}$$

- $P(\vec{x}|R = 1, \vec{q})$ and $P(\vec{x}|R = 0, \vec{q})$: probability that if a relevant or nonrelevant document is retrieved, then that document's representation is $\vec{x}$
- Use statistics about the document collection to estimate these probabilities

# IIR 12: Language models



| $w$ | $P(w|q_1)$ | $w$ | $P(w|q_1)$ |
|------|--------|-------|--------|
| STOP | 0.2 | toad | 0.01 |
| the | 0.2 | said | 0.03 |
| a | 0.1 | likes | 0.02 |
| frog | 0.01 | that | 0.04 |
| | | . . . | . . . |

This is a one-state probabilistic finite-state automaton – a unigram language model – and the state emission distribution for its one state $q_1$.

STOP is not a word, but a special symbol indicating that the automaton stops.

"frog said that toad likes frog STOP"

$P(\text{string}) = 0.01 \cdot 0.03 \cdot 0.04 \cdot 0.01 \cdot 0.02 \cdot 0.01 \cdot 0.2 = 0.0000000000048$

# IIR 13: Text classification & Naive Bayes

- Text classification = assigning documents automatically to predefined classes
- Examples:
  - Language (English vs. French)
  - Adult content
  - Region
  - Fake news
  - Jeopardy challenge
  - Multiple choice questions

# IIR 14: Vector classification

# IIR 15: Learning to rank

- How to rank documents/pages when they are described by multiple pieces of information (features).
- Use machine learning to quantify the importance of each feature.
- Use this for your project if you are a grad student!

# IIR 16: Flat clustering

# IIR 17: Hierarchical clustering

http://news.google.com



**Barack Obama, Vladimir Putin hold first in-person talk since start of Ukraine ...**
Canada.com - 6 hours ago
A video screen shows U.S. President Barack Obama and Russian President Vladimir Putin during the Ouistreham commemorate the 70th anniversary of the Allied invasion at Normandy, in Ouistreham, France, June ...

**Obama and Vladimir Putin dine separately with French President Francois ...**
New York Daily News - 10 hours ago
French President Francois Hollande engaged in some deft dinnertime diplomacy Thursday night - hosting separate meals so Pr Vladimir Putin wouldn't have to break bread together. Hollande and Obama had dinner at ...



**Topless feminist stabs wax Putin in France**
The Local.fr - Jun 5, 2014
The same day President Vladimir Putin was to arrive in France for D-Day anniversary events, radical feminist p leader's statue in a Paris wax museum.

**Vladimir Putin meets Ukraine president-elect Petro Poroshenko at D-Day ...**
Telegraph.co.uk - 9 hours ago
Russian President Vladimir Putin and Ukraine's president-elect Petro Poroshenko discussed a ceasefire and other possible ste countries in a brief but potentially significant meeting in France on Friday, French ...
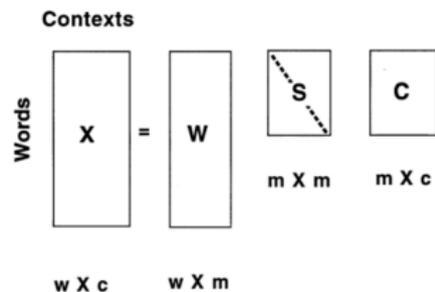


**Vote Possible to Get Stalingrad Name Back: Vladimir Putin**
NDTV - 2 hours ago
Russian President Vladimir Putin speaks to the media at Benouville castle, Friday, June 6, 2014, where he arriv landings.
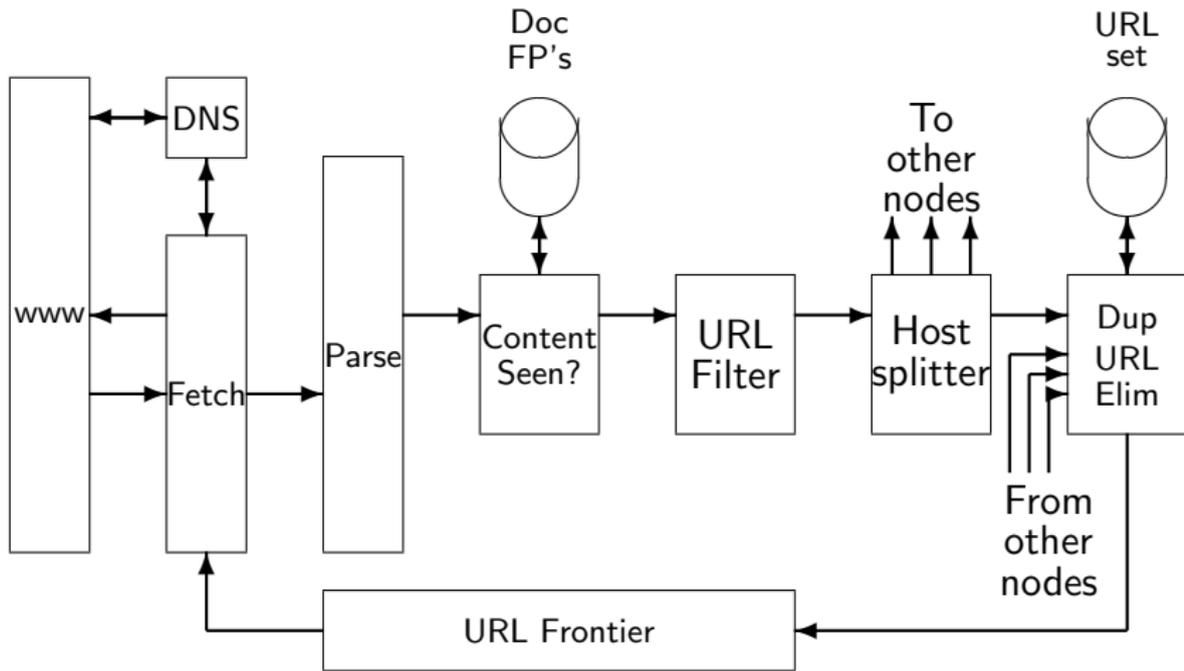
# IIR 18: Latent Semantic Indexing

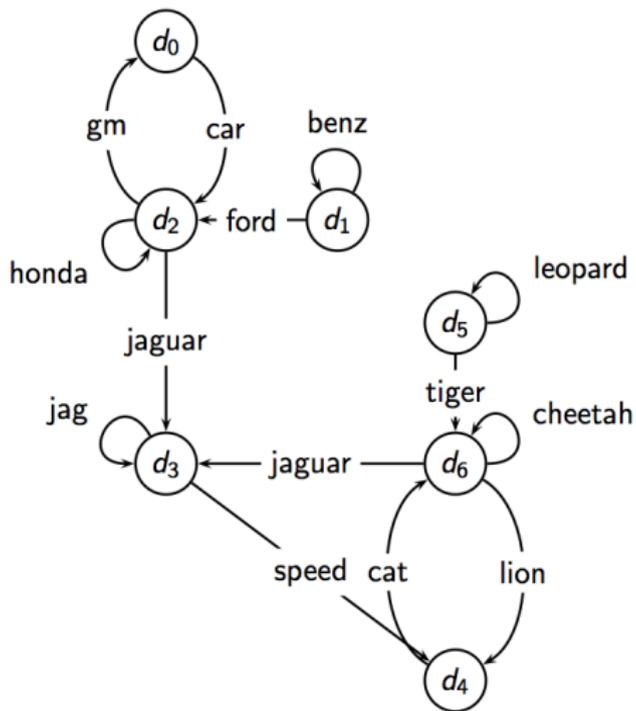"You shall know a word by the company it keeps."

## IIR 19: The web and its challenges

- Unusual and diverse documents
- Unusual and diverse users and information needs
- Beyond terms and text: exploit link analysis, user data
- How do web search engines work?
- How can we make them better?

# IIR 20: Crawling

# IIR 21: Link analysis / PageRank

## Take-away

- Why you should take this course
- Admin issues
- Boolean Retrieval: Design and data structures of a simple information retrieval system
- What topics will be covered in this class?