

# Introduction to Information Retrieval

(Manning, Raghavan, Schutze)

## Chapter 1

### Boolean retrieval

# Information Retrieval: IR

---

- Finding material (usually document) of an unstructured nature (usually text) that satisfies an information need from within large collections
- Started in the 50' s. SIGIR (80), TREC (92)
- The field of IR also covers supporting users in browsing or filtering document collections or further processing a set of retrieved documents
  - clustering
  - classification
- Scale: from web search to personal information retrieval

# How good are the retrieved docs?

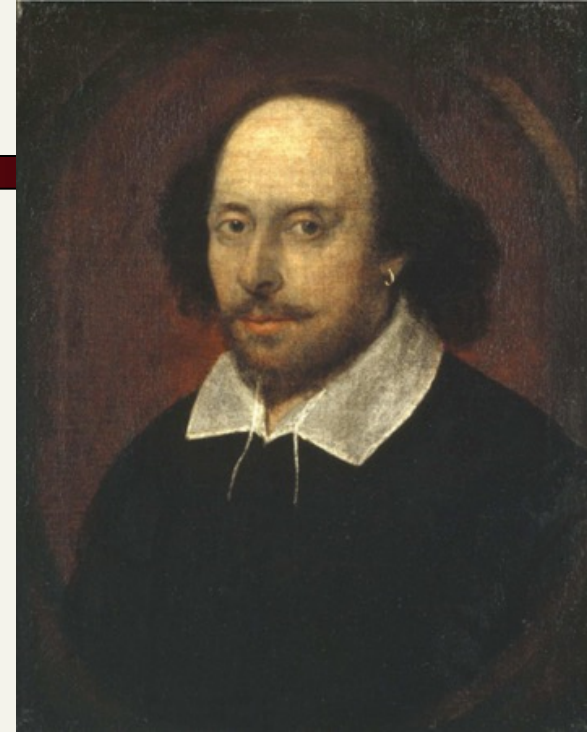
---

- Precision : Fraction of retrieved docs that are relevant to user's information need
- Recall : Fraction of relevant docs in collection that are retrieved
- More precise definitions and measurements to follow in later lectures

# Boolean retrieval

---

- Queries are Boolean expressions
  - e.g., *Brutus AND Caesar*
  - Shakespeare's Collected Works
  - Which plays of Shakespeare contain the words *Brutus AND Caesar*?
- The search engine returns all documents satisfying the Boolean expression.
  - Does Google use the Boolean model?
- <http://www.rhymezone.com/shakespeare/>



# Example

---

- 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 corpora)
  - “Not Calpurnia” is non-trivial
  - Other operations (e.g., find the word *Romans* near *countrymen*) not feasible
  - Ranked retrieval (best documents to return)

# Term-document incidence

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

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

1 if play contains  
word, 0 otherwise

# Incidence vectors

---

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

# Answers to query

---

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

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

## ■ Hamlet, Act III, Scene ii

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



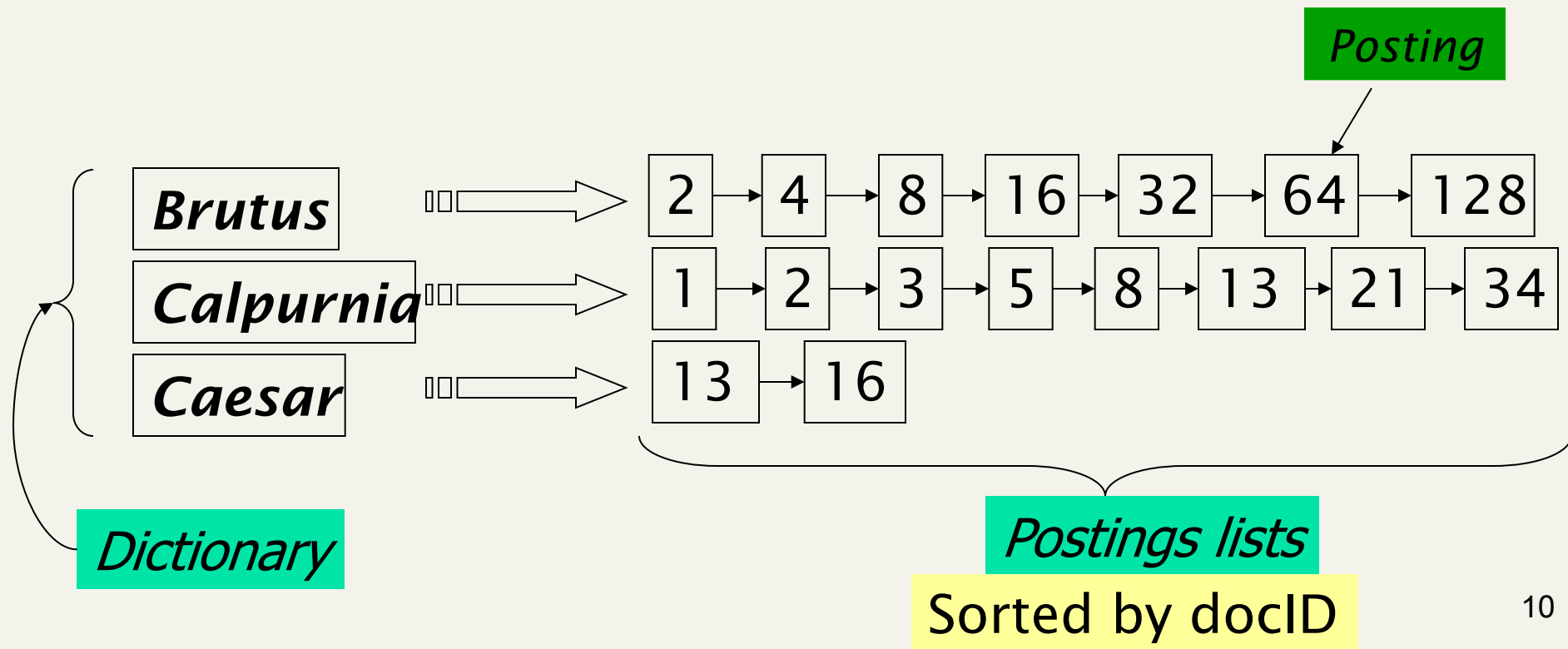
# Big collections: cannot build the matrix

---

- Consider  $N = 1\text{M}$  documents, each with about 1K terms
- Avg 6 bytes/term including spaces/punctuation
  - Size of document collection is about 6GB
- Say there are  $m = 500\text{K}$  distinct terms among these.
- 500K x 1M matrix has half-a-trillion 0's and 1's.
- But it has no more than one billion 1's.
  - matrix is extremely sparse.
- What's a better representation?
  - We only record the 1s

# Inverted index

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



# Inverted index construction

Documents to be indexed.



Friends, Romans, countrymen.  
⋮

Tokenizer

Token stream.

Friends

Romans

Countrymen

Linguistic modules

Modified tokens.

friend

roman

countryman

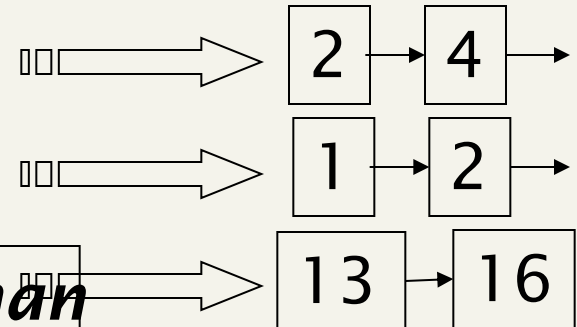
Indexer

Inverted index.

**friend**

**roman**

**countryman**





# Indexer steps

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

Doc 1

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

Doc 2

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



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



## ■ Sort by terms.

**Core indexing step.**

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



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



- Multiple term entries in a single document are merged.
- Frequency information is added.

Why frequency?

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



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

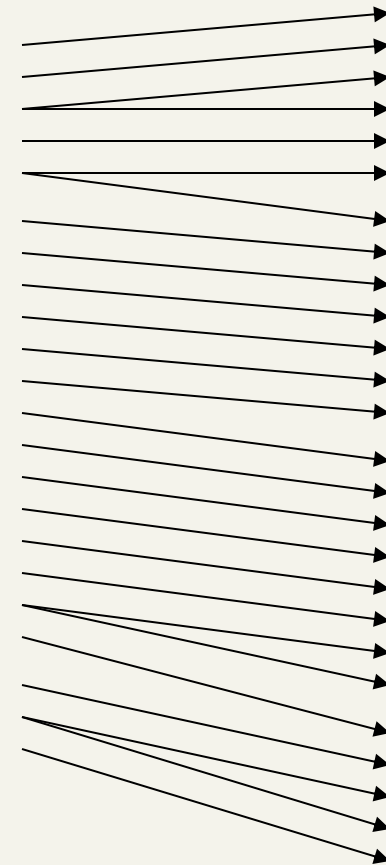


- The result is split into a *Dictionary* file and a *Postings* file.

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



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



Doc #	Freq
2	1
2	1
1	1
2	1
1	1
1	1
2	2
1	1
1	1
2	1
1	2
1	1
2	1
1	1
2	1
2	1
1	1
2	1
2	1
2	1
1	1
2	1
2	1
1	1
2	1
2	1

## ■ Where do we pay in storage?

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

Terms

Pointers

Document  
frequency

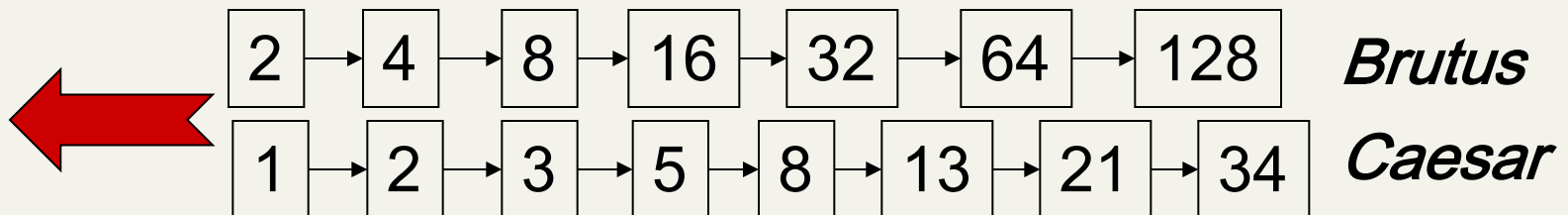
Collection  
frequency

Term  
frequency



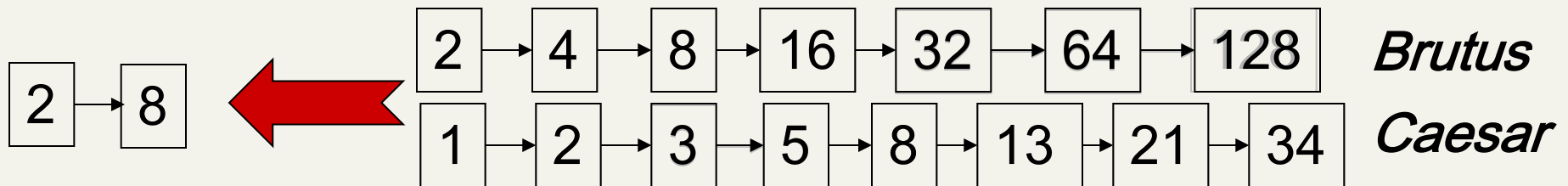
# Boolean query processing: AND

- Consider processing the query:  
*Brutus AND Caesar*
  - Locate *Brutus* in the Dictionary;
    - Retrieve its postings.
  - Locate *Caesar* in the Dictionary;
    - Retrieve its postings.
  - “Merge” the two postings:



# 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 sorted by docID.

# Boolean queries: Exact match

---

- The Boolean Retrieval model is being able to ask a query that is a Boolean expression:
  - Boolean Queries are queries using *AND*, *OR* and *NOT* to join query terms
    - Views each document as a set of words
    - Is precise: document matches condition or not.
- Primary commercial retrieval tool for 3 decades.
- Professional searchers (e.g., lawyers) still like Boolean queries:
  - You know exactly what you're getting.
- Many search systems you use are Boolean
  - Email, Intranet etc.

# Example: WestLaw

<http://www.westlaw.com/>

- Commercially successful Boolean retrieval
- Largest commercial (paying subscribers) legal search service (started 1975; ranking added 1992)
- Tens of terabytes of data; 700,000 users
- Majority of users *still* use boolean queries
- Example query:
  - What is the statute of limitations in cases involving the federal tort claims act?
  - **LIMIT! /3 STATUTE ACTION /S FEDERAL /2 TORT /3 CLAIM**
- /3 = within 3 words, /S = in same sentence

# Example: WestLaw

<http://www.westlaw.com/>

- Another example query:
  - Requirements for disabled people to be able to access a workplace
  - `disabl! /p access! /s work-site work-place (employment /3 place`
- Note that SPACE is disjunction, not conjunction!
- Long, precise queries; proximity operators; incrementally developed; not like web search
- Professional searchers often like Boolean search:
  - Precision, transparency and control
- But that doesn't mean they actually work better....

# Boolean queries: more general merges

---

- Exercise: Adapt the merge for the queries:
  - (a) *Brutus AND NOT Caesar*
  - (b) *Brutus OR NOT Caesar*

Can we still run through the merge in time  $O(x+y)$ ?  
What can we achieve?

# Merging

---

What about an arbitrary Boolean formula?

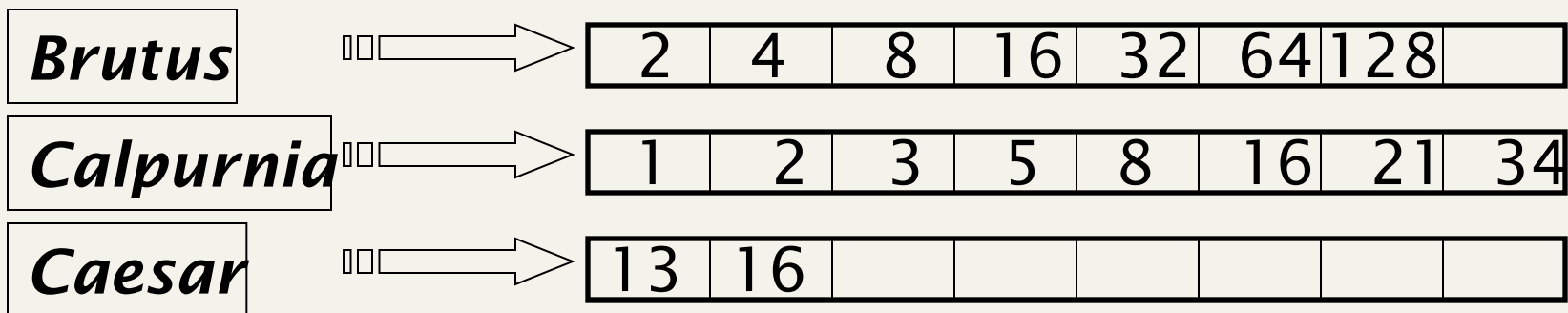
*(Brutus OR Caesar) AND NOT*

*(Antony OR Cleopatra)*

- Can we always merge in “linear” time?
  - Linear in what?
- Can we do better?

# Query optimization

- What is the best order for query processing?
- Consider a query that is an *AND* of  $t$  terms.
- For each of the  $t$  terms, get its postings, then *AND* them together.



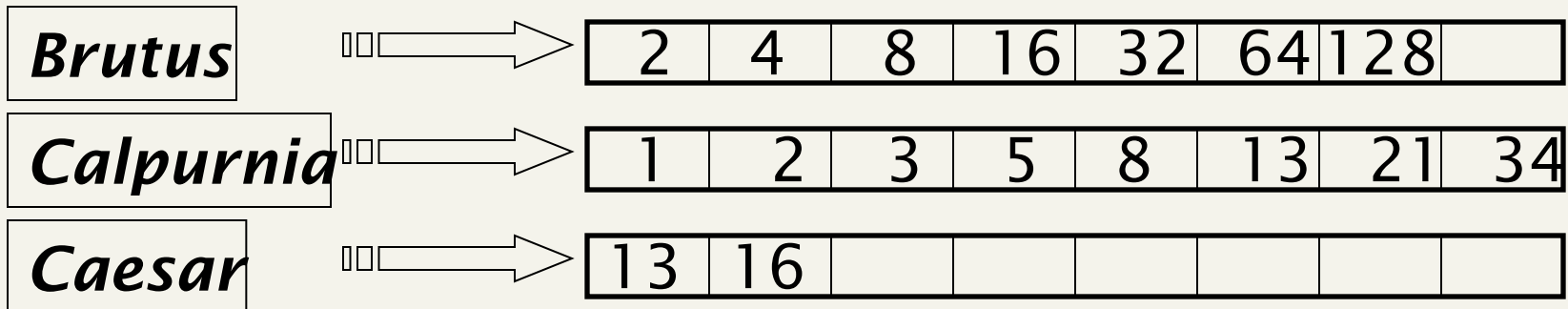
Query: **Brutus AND Calpurnia AND Caesar**



# Query optimization example

- Process in order of increasing freq:
  - *start with smallest set, then keep cutting further.*

This is why we kept  
freq in dictionary



Execute the query as  $(Caesar \text{ AND } Brutus) \text{ AND } Calpurnia$ .

# More general optimization

---

- e.g., (*madding OR crowd*) AND (*ignoble OR strife*)
- Get freq' s for all terms.
- Estimate the size of each *OR* by the sum of its freq' s (conservative).
- Process in increasing order of *OR* sizes.

# Exercise

- Recommend a query processing order for

*(tangerine OR trees) AND  
(marmalade OR skies) AND  
(kaleidoscope OR eyes)*

Kaleidoscope OR eyes (300,321)

Tangerine OR trees (363,465)

Marmalade OR skies (379,571)

Term	Freq
eyes	213312
kaleidoscope	87009
marmalade	107913
skies	271658
tangerine	46653
trees	316812

# Coverage

---

- 1: Introduction
- 2: Term vocabulary and postings lists
- 3: Dictionaries (briefly)
- 4: Index construction (briefly)
- 5: Index compression (briefly)
- 6: Term weighting and vector space model**
- 7: Computing scores (briefly)
- 8: Evaluation and result summaries**
- 19: Web search basics
- 21: Link analysis**