**CalPoly**Pomona

# Indexes

Lecture 5

1

# Indexes

- *Indexes* are data structures designed to make search faster
- Text search has unique requirements, which leads to unique data structures
- Most common data structure is *inverted index*
  - "inverted" because documents are associated with words, rather than words with documents
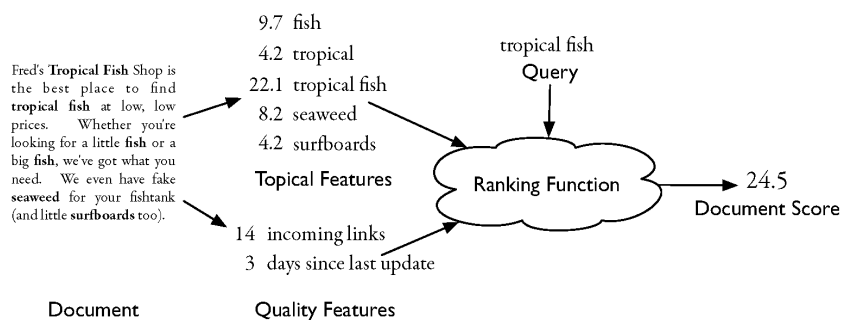    - similar to a *concordance*

2

# Indexes and Ranking

- Indexes are designed to support *search*
  - faster response time, supports updates
- Text search engines use a particular form of search: *ranking*
  - documents are retrieved in sorted order according to a score computing using the document representation, the query, and a *ranking algorithm*
- What is a reasonable abstract model for ranking?
  - enables discussion of indexes without details of retrieval model

3

# Abstract Model of Ranking

Fred's **Tropical Fish** Shop is the best place to find **tropical fish** at low, low prices. Whether you're looking for a little **fish** or a big **fish**, we've got what you need. We even have fake **seaweed** for your fishtank (and little **surfboards** too).

Document

9.7  fish
4.2  tropical
22.1  tropical fish
8.2  seaweed
4.2  surfboards

Topical Features

14  incoming links
3  days since last update

Quality Features

tropical fish
Query

Ranking Function

24.5
Document Score

4

2

# More Concrete Model

$$R(Q, D) = \sum_i g_i(Q) f_i(D)$$

$f_i$ is a document feature function
$g_i$ is a query feature function



Fred's **Tropical Fish** Shop is the best place to find **tropical fish** at low, low prices. Whether you're looking for a little **fish** or a big **fish**, we've got what you need. We even have fake **seaweed** for your fishtank (and little **surfboards** too).

Document

**f_i**

| | Topical Features |
|---|---|
| 9.7 | fish |
| 4.2 | tropical |
| 22.1 | tropical fish |
| 8.2 | seaweed |
| 4.2 | surfboards |

Quality Features

| 14 | incoming links |
|---|---|
| 3 | update count |

**g_i**

| Topical Features | |
|---|---|
| fish | 5.2 |
| tropical | 3.4 |
| tropical fish | 9.9 |
| chichlids | 1.2 |
| barbs | 0.7 |

tropical fish
Query

Quality Features

| incoming links | 1.2 |
|---|---|
| update count | 0.9 |

303.01
Document Score

# Inverted Index

- Each index term is associated with an *inverted list*
  - Contains lists of documents, or lists of word occurrences in documents, and other information
  - Each entry is called a *posting*
  - The part of the posting that refers to a specific document or location is called a *pointer*
  - Each document in the collection is given a unique number
  - Lists are usually *document-ordered* (sorted by document number)
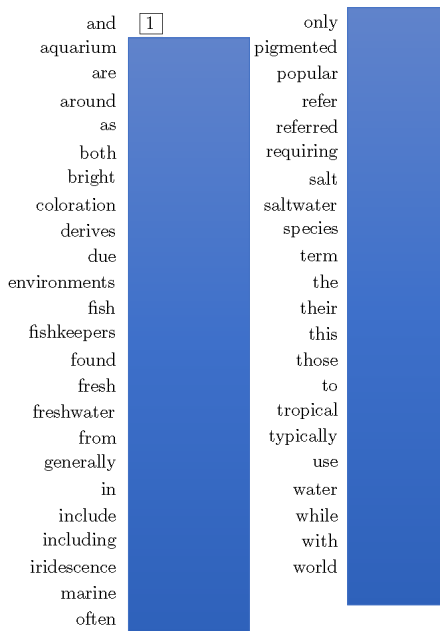
# Example "Collection"

$S_1$   Tropical fish include fish found in tropical environments around the world, including both freshwater and salt water species.

$S_2$   Fishkeepers often use the term tropical fish to refer only those requiring fresh water, with saltwater tropical fish referred to as marine fish.

$S_3$   Tropical fish are popular aquarium fish, due to their often bright coloration.

$S_4$   In freshwater fish, this coloration typically derives from iridescence, while salt water fish are generally pigmented.

Four sentences from the Wikipedia entry for *tropical fish*

7

---

## Simple Inverted Index

| | [1] | |
|---|---|---|
| and | | only |
| aquarium | | pigmented |
| are | | popular |
| around | | refer |
| as | | referred |
| both | | requiring |
| bright | | salt |
| coloration | | saltwater |
| derives | | species |
| due | | term |
| environments | | the |
| fish | | their |
| fishkeepers | | this |
| found | | those |
| fresh | | to |
| freshwater | | tropical |
| from | | typically |
| generally | | use |
| in | | water |
| include | | while |
| including | | with |
| iridescence | | world |
| marine | | |
| often | | |

8

4

# Simple Inverted Index

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| and | 1 | | | | only | 2 | | |
| aquarium | 3 | | | | pigmented | 4 | | |
| are | 3 | 4 | | | popular | 3 | | |
| around | 1 | | | | refer | 2 | | |
| as | 2 | | | | referred | 2 | | |
| both | 1 | | | | requiring | 2 | | |
| bright | 3 | | | | salt | 1 | 4 | |
| coloration | 3 | 4 | | | saltwater | 2 | | |
| derives | 4 | | | | species | 1 | | |
| due | 3 | | | | term | 2 | | |
| environments | 1 | | | | the | 1 | 2 | |
| fish | 1 | 2 | 3 | 4 | their | 3 | | |
| fishkeepers | 2 | | | | this | 4 | | |
| found | 1 | | | | those | 2 | | |
| fresh | 2 | | | | to | 2 | 3 | |
| freshwater | 1 | 4 | | | tropical | 1 | 2 | 3 |
| from | 4 | | | | typically | 4 | | |
| generally | 4 | | | | use | 2 | | |
| in | 1 | 4 | | | water | 1 | 2 | 4 |
| include | 1 | | | | while | 4 | | |
| including | 1 | | | | with | 2 | | |
| iridescence | 4 | | | | world | 1 | | |
| marine | 2 | | | | | | | |
| often | 2 | 3 | | | | | | |

# Inverted Index with counts

- supports better ranking algorithms

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| and | 1:1 | | | | only | 2:1 | | |
| aquarium | 3:1 | | | | pigmented | 4:1 | | |
| are | 3:1 | 4:1 | | | popular | 3:1 | | |
| around | 1:1 | | | | refer | 2:1 | | |
| as | 2:1 | | | | referred | 2:1 | | |
| both | 1:1 | | | | requiring | 2:1 | | |
| bright | 3:1 | | | | salt | 1:1 | 4:1 | |
| coloration | 3:1 | 4:1 | | | saltwater | 2:1 | | |
| derives | 4:1 | | | | species | 1:1 | | |
| due | 3:1 | | | | term | 2:1 | | |
| environments | 1:1 | | | | the | 1:1 | 2:1 | |
| fish | 1:2 | 2:3 | 3:2 | 4:2 | their | 3:1 | | |
| fishkeepers | 2:1 | | | | this | 4:1 | | |
| found | 1:1 | | | | those | 2:1 | | |
| fresh | 2:1 | | | | to | 2:2 | 3:1 | |
| freshwater | 1:1 | 4:1 | | | tropical | 1:2 | 2:2 | 3:1 |
| from | 4:1 | | | | typically | 4:1 | | |
| generally | 4:1 | | | | use | 2:1 | | |
| in | 1:1 | 4:1 | | | water | 1:1 | 2:1 | 4:1 |
| include | 1:1 | | | | while | 4:1 | | |
| including | 1:1 | | | | with | 2:1 | | |
| iridescence | 4:1 | | | | world | 1:1 | | |
| marine | 2:1 | | | | | | | |
| often | 2:1 | 3:1 | | | | | | |

## Inverted Index with positions

- supports proximity matches

| Term | Postings | | | | |
|------|------|------|------|------|------|
| and | 1,15 | | | | |
| aquarium | 3,5 | | | | |
| are | 3,3 | 4,14 | | | |
| around | 1,9 | | | | |
| as | 2,21 | | | | |
| both | 1,13 | | | | |
| bright | 3,11 | | | | |
| coloration | 3,12 | 4,5 | | | |
| derives | 4,7 | | | | |
| due | 3,7 | | | | |
| environments | 1,8 | | | | |
| fish | 1,2 | 1,4 | 2,7 | 2,18 | 2,23 |
| | 3,2 | 3,6 | 4,3 | | |
| | 4,13 | | | | |
| fishkeepers | 2,1 | | | | |
| found | 1,5 | | | | |
| fresh | 2,13 | | | | |
| freshwater | 1,14 | 4,2 | | | |
| from | 4,8 | | | | |
| generally | 4,15 | | | | |
| in | 1,6 | 4,1 | | | |
| include | 1,3 | | | | |
| including | 1,12 | | | | |
| iridescence | 4,9 | | | | |

| Term | Postings | | | | |
|------|------|------|------|------|------|
| marine | 2,22 | | | | |
| often | 2,2 | 3,10 | | | |
| only | 2,10 | | | | |
| pigmented | 4,16 | | | | |
| popular | 3,4 | | | | |
| refer | 2,9 | | | | |
| referred | 2,19 | | | | |
| requiring | 2,12 | | | | |
| salt | 1,16 | 4,11 | | | |
| saltwater | 2,16 | | | | |
| species | 1,18 | | | | |
| term | 2,5 | | | | |
| the | 1,10 | 2,4 | | | |
| their | 3,9 | | | | |
| this | 4,4 | | | | |
| those | 2,11 | | | | |
| to | 2,8 | 2,20 | 3,8 | | |
| tropical | 1,1 | 1,7 | 2,6 | 2,17 | 3,1 |
| typically | 4,6 | | | | |
| use | 2,3 | | | | |
| water | 1,17 | 2,14 | 4,12 | | |
| while | 4,10 | | | | |
| with | 2,15 | | | | |
| world | 1,11 | | | | |

CS 4250 – Web Search and Recommender Systems

11

# Proximity Matches

- Matching phrases or words within a window
  - e.g., "tropical fish", or "find tropical within 5 words of fish"
- Word positions in inverted lists make these types of query features efficient
  - e.g.,

| tropical | 1,1 | | 1,7 | 2,6 | 2,17 | | 3,1 | | |
|----------|-----|-----|-----|-----|------|------|-----|-----|------|
| fish | 1,2 | 1,4 | | 2,7 | 2,18 | 2,23 | 3,2 | 3,6 | 4,3 | 4,13 |

CS 4250 – Web Search and Recommender Systems

12

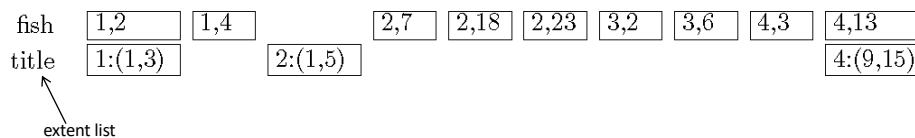# Fields and Extents

- Document structure is useful in search
  - *field* restrictions
    - e.g., date, from:, etc.
  - some fields more important
    - e.g., title
- Options:
  - separate inverted lists for each field type
  - add information about fields to postings
  - use *extent lists*

13

# Extent Lists

- An *extent* is a contiguous region of a document
  - represent extents using word positions
  - inverted list records all extents for a given field type
  - e.g.,

| fish | 1,2 | 1,4 | | 2,7 | 2,18 | 2,23 | 3,2 | 3,6 | 4,3 | 4,13 |
| title | 1:(1,3) | | 2:(1,5) | | | | | | | 4:(9,15) |

extent list

14

# Precomputed Scores

- Precomputed scores in inverted list
    - e.g., list for "fish" [(1:3.6), (3:2.2)], where 3.6 is total feature value for document 1
    - Score could be based on many different attributes (frequency, title occurrence, etc.)
    - improves speed but reduces flexibility
    - What is lost?
        - number of terms
        - term position
        - ...

15

# Compression

- Inverted lists are very large
    - e.g., 25-50% of collection for TREC collections using Indri search engine
    - Much higher if n-grams are indexed
- Compression of indexes saves disk and/or memory space
    - Typically have to decompress lists to use them
    - Best compression techniques have good *compression ratios* and are easy to decompress

16

8

# Query Processing

- Document-at-a-time
    - Calculates complete scores for documents by processing all term lists, one document at a time
- Term-at-a-time
    - Accumulates scores for documents by processing term lists one at a time
- Both approaches have optimization techniques that significantly reduce time required to generate scores

# Document-At-A-Time

Query: *salt water tropical*

| salt | 1:1 | | | 4:1 |
| water | 1:1 | 2:1 | | 4:1 |
| tropical | 1:2 | 2:2 | 3:1 | |
| **score** | 1:4 | 2:3 | 3:1 | 4:2 |

# Term-At-A-Time

Query: *salt water tropical*

| | | |
|---|---|---|
| salt | 1:1 | 4:1 |
| partial scores | 1:1 | 4:1 |

| | | | |
|---|---|---|---|
| old partial scores | 1:1 | | 4:1 |
| water | 1:1 | 2:1 | 4:1 |
| new partial scores | 1:2 | 2:1 | 4:2 |

| | | | | |
|---|---|---|---|---|
| old partial scores | 1:2 | 2:1 | | 4:2 |
| tropical | 1:2 | 2:2 | 3:1 | |
| final scores | 1:4 | 2:3 | 2:2 | 4:2 |

CS 4250 – Web Search and Recommender Systems

19

# Caching

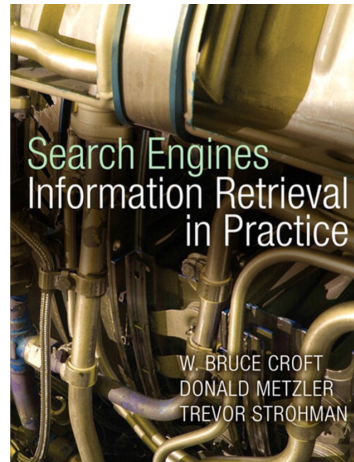- Query distributions similar to Zipf
  - About ½ each day are unique, but some are very popular
- Caching can significantly improve effectiveness
  - E.g. Cache popular query results
- Cache must be refreshed to prevent stale data

CS 4250 – Web Search and Recommender Systems

20

# Reading

• Chapter 5

Search Engines
Information Retrieval
in Practice

W. BRUCE CROFT
DONALD METZLER
TREVOR STROHMAN

CS 4250 – Web Search and Recommender Systems

21