Applications of Indexing

Lecture 16



Applications of Indexing

Overview

- We've talked about how indexes can improve performance inside a DBMS
- But these same data structures/ techniques can be useful on other systems as well!
- Let's look at a couple examples :)



Outline

- 1. Problem #1: What is **Full-Text Search**?
 - Can it be done naïvely? (spoiler: not for Google!)
- 2. To the Rescue: The Inverted Index
 - Design a relational index
 - Advanced Issues
 - Research Application: Cognitive Modeling
- 3. The R-tree
 - Overview
 - Research Application: Optimization



Applications of Indexing

What is Full-Text Search?

<u>Given</u>

- A set of "documents"
 - Each containing "words"

Problem

- Find the "best" document(s) that satisfy a query
 - Simplest: set of words

<u>Requirements</u>

- Fast & scalable
- Relevant results
- Expressive queries
- Up-to-date



Applications of Indexing

Example: Web Search



Other Examples





search

northeastern

advanced search: by author, subreddit...

subreddits

Northeastern University

subscriber r/northeastern 190 subscribers, a community for 7 years p search within r/northeastern

Northeastern University (NU/NEU)

subscribe r/NEU 5,348 subscribers, a community for 8 years News and discussion of interest to students, faculty, employees, and neighbors of **Northeastern** University in Boston, MA. ρ search within r/NEU



Applications of Indexing

Fast & Scalable: # of Documents

Google	> 130 Trillion Pages
facebook	~ 2 Billion Users
amazon	~ 400 Million Products



Applications of Indexing

Fast & Scalable: Queries/sec





Applications of Indexing

Big-O Review

- What does it mean for an algorithm to scale in linear time (i.e. O(*n*))?
- Describe a linear-time algorithm for fulltext search of keywords in webpages
 - Find all documents that contain "northeastern"



~10.4M blu ray ~ 7.75 miles



Assume simple query:

northeastern

- Single repo of all pages ~ 15 x Burj Khalifa
 130T * 250 w/page * 8 bytes/w ~ 260 PB
- Require 1s response time
 - 8M * 3GHz 64-bit CPU (assume 1 cycle/w)
 - (8M * 63K) CPUs * 35W/CPU * \$0.14/kWH

~ \$21.6T/year ~ 1.25 x US GDP

~ \$0.7M/sec

~67 CPU/person



Applications of Indexing

Indexing

- Improve search speed at the cost of extra...
 - **Memory** for data structure(s)
 - **Time** to update the data structure(s)
- Backbone of search engines, databases, graphics/game engines, simulation software, ...



Consider a Physical Textbook

 Find all pages of the book that contain the word "husky"





Official Publication of the AMERICAN KENNEL CLUB More than 2 million copies sold



Analyze a Physical (Inverted) Index

1. Find the term

- Alphabetical search
 - More generally: binary
- Time: O(log₂terms)

2. Enumerate the pages

- Linear scan
 - Page 1, 2, ..., *k*
- Time: O(pages)

Space: O(terms + pages)



Applications of Indexing

Linear vs Logarithmic (1)

	N=10	N=100	N=1000
Logarithmic	1s	2s	3s
Linear	1s	10s	100s



Linear vs Logarithmic (2)





Linear vs Logarithmic (3)





Applications of Indexing

Inverted Index by Example

Given documents { D_1 , D_2 , D_3 }:

- $-D_1 =$ "it is what it is"
- $-D_2 =$ "what is it"
- $-D_3 =$ "it is a banana"

Inverted Index:





Applications of Indexing

Design a Relational Inverted Index

Develop a set of table(s) and index(es) that support efficient construction and querying of an inverted index

Documents		Contents	
Docld	DocPath	Word	<u>Docld</u>
1	/path/to/doc	Foo	1

SELECT DocId FROM Contents WHERE Word=? ORDER BY DocId



Applications of Indexing

Advanced Issues

- More expressive query semantics
 - Multi-word
 - Locality: ["what is it"] vs. ["what it is"] vs. ["what", "is", "it"]
- Ranked results
 - Document-ranking algorithm (e.g. PageRank)
 - Efficient ranked retrieval
- Dynamics
 - Document addition/removal/modification
 - Rank
 - Document changes
 - Integration of real-time variables (e.g. location)



Applications of Indexing

Research: Cognitive Modeling

- Semantic Memory is a human's long-term store of facts about the world, independent of the context in which they were originally learned
- The ACT-R (<u>http://act-r.psy.cmu.edu</u>) model of semantic memory has been successful at explaining a variety of psychological phenomena (e.g. retrieval bias, forgetting)
- The model does not scale to large memory sizes, which hampers complex experiments



Applications of Indexing

Memory Representation



- Document = Node
- Word = edge
- Office hours: ask me about the number inside each node :)

Example cue: last(obama), spouse(X)



Applications of Indexing

Scale Fail [AFRL '09]

Retrieval Latency: Chunks in DM x Retrieval Constraints x Type of DM (Error Bars: 95% Confidence Interval)



Applications of Indexing

Example

Semantic Objects: Features





Applications of Indexing

Inverted Index



Index Statistics



Top-1 Retrieval

Inverted Index



Some Results [ICCM '10]

Inverted index (via SQLite) + new approach was **fast** and **scaled**!

>30x faster than off-the-shelf database (on >3x data)!



Applications: AI + Inverted Index

Learned Navigation







Task Learning



Applications of Indexing

Research: Optimization

Packing. Fit n circles of radius r in a square of side-length s without overlap (non-convex, NP-hard, ∞ solutions). Used in making codes, physical packing, computer-assisted origami.







Applications of Indexing

Another Index: R-tree

"Group nearby objects and represent them with their **minimum bounding rectangle** in the next higher level of the tree... Since all objects lie within this bounding rectangle, a query that does not intersect the bounding rectangle also cannot intersect any of the contained objects. At the leaf level, each rectangle describes a single object; at higher levels the aggregation of an increasing number of objects." -- Wikipedia





Core Idea: reduce the constraints via relative distance



Applications of Indexing

Large-Scale Evaluation [BICA '14]





Applications of Indexing

Takeaways

- A common approach to large-scale search is indexing: using data structure(s) to improve access speed
- An inverted index is commonly used for full-text search (even in situations that might not look like it)
 - Inverted indexes are fast, scalable, and straight-forward to implement
- An R-tree is commonly used for spatial queries over objects in 2/3D space (e.g. what is within X miles of Y? are A and B colliding?)
- <u>Know your indexes/data structures</u> careful problem analysis and algorithm development can often beat generic approaches
 - Even if you don't use a DBMS, DBMS methods can be very useful in a variety of applications!

