# Search Engines
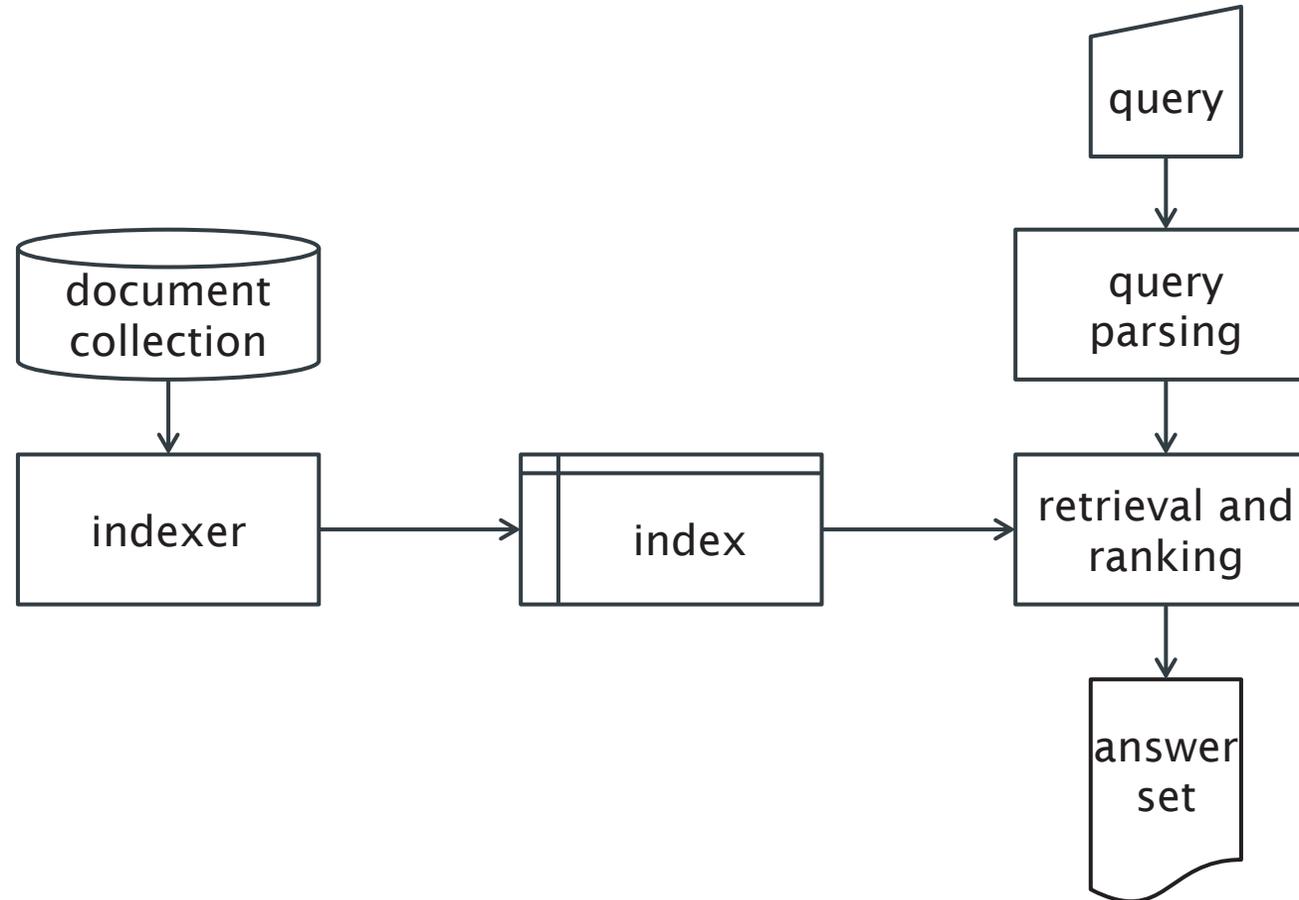
COMP3227 Web Architecture & Hypertext Technologies
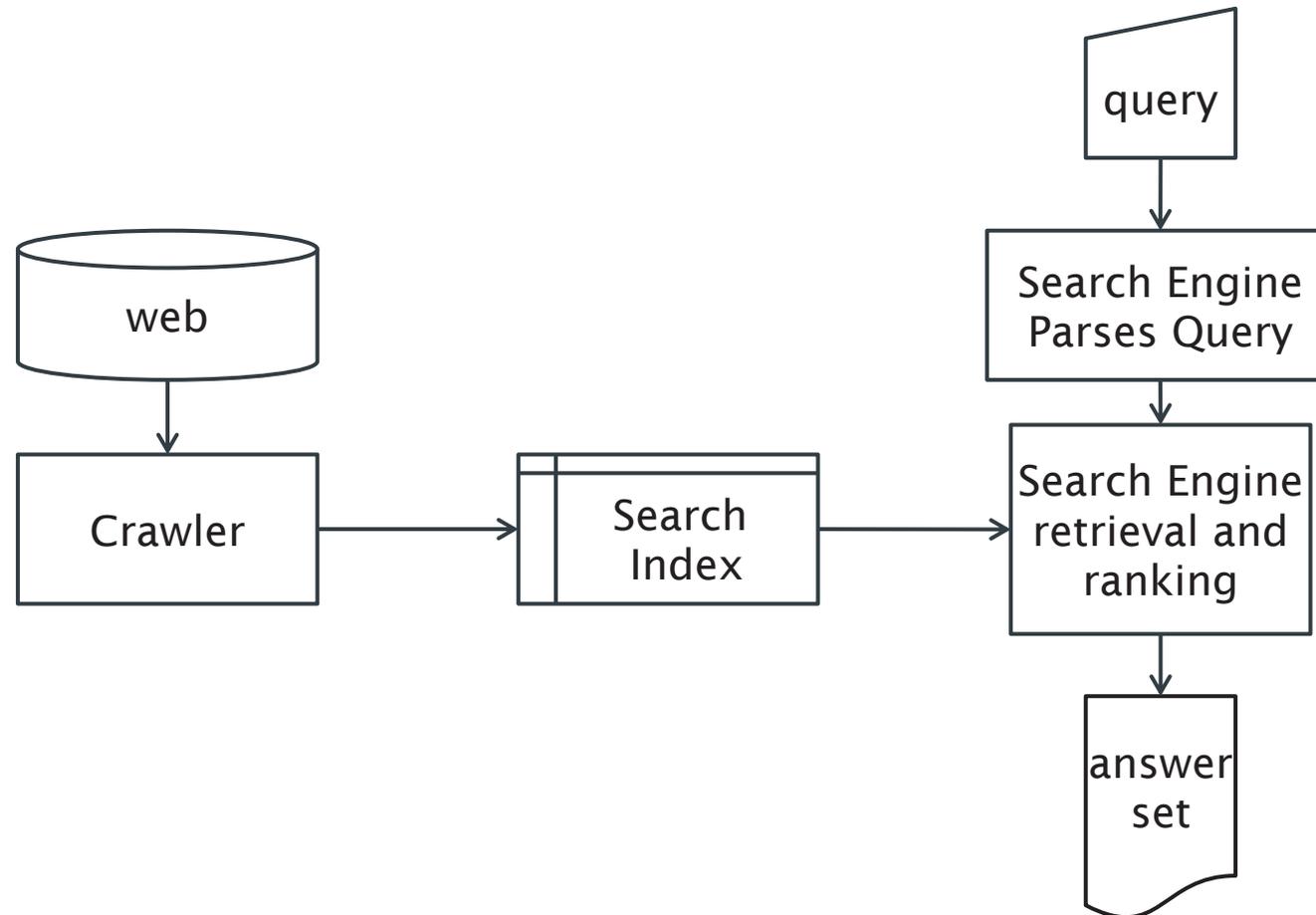
Dr Heather Packer – hp3@ecs.soton.ac.uk

# Information Retrieval

- The primary goal of an information retrieval system is to retrieve all the documents that are relevant to a user query while retrieving as few non-relevant documents as possible

  – An **information need** is a topic which a user desires to know more about

  – A **query** is what the user conveys to the computer in an attempt to communicate their information need

  – A document is **relevant** if it is one that the user perceives as containing information of value with respect to their personal information need

# High-Level System Architecture

# Search Engine High-Level System Architecture

# Specific Board and Vague Queries

| | Specific | Broad | Vague |
|---|---|---|---|
| I know specifically what I'm looking for | ✔ | ✘ | ✘ |
| I know where to look for what I'm looking for | ✔ | ✔ | ✘ |
| Example | Looking for a specific gene in the Gene Database | Looking for the manager of HR, in the company directory | Google |

Search Engines

# Search Engines

# Search Engines

- Search engines are a service

- They allow users to search for content using **keywords**

- A query returns a **ranked** set of results

- They **DO NOT** access the web directly

- They **USE** huge databases

# Web Crawler, Spider or bot

- An algorithm that systematically browses the web

- A basic algorithm

    1) Start at a webpage

    2) Follow the hyperlinks that webpage points to

    3) Then follows the links those webpages point to

- Each page it visits it collects metadata about it

- Stores a file for each resource, with its metadata in a search index

- Crawlers consume resources

- Can visit sites without approval

# **Web Crawler** - robots.txt

- Block all web crawlers

User-agent: *

Disallow: /

- Allow all web crawlers

User-agent: *

Disallow:

- Block a specific web crawler from a specific folder

User-agent: Googlebot

Disallow: /example-subfolder/

Disallow: /index.html

# Web Crawler Policies

The behaviour of a web crawler is based on:

1. Selection Policy

2. Re-visit Policy

3. Politeness Policy

4. Parallelisation Policy

# **Web Crawler** Selection Policy

- Search engines only index part of the web

- It's important to download the most relevant pages

- A selection policy states which pages to index

- Strategies:

  – Breadth first

  – Back link count

  – PageRank

- Only request pages that have searchable content (HTML, PDF etc)

# Web Crawler Re-visit Policy

- Its worth revisiting web pages because they change over time

- An ideal search engine would have the most up-to-date version of every page in its index

- Strategies

  - Re-visit all pages equally frequent

  - Prioritise pages that change often (but not too often!)

- May take page quality into account

# **Web Crawler** Politeness Policy

- Issues of schedule and load when large collections of pages are accessed

- Strategies

  - Do not make parallel calls to the same server

  - Spread out requests

  - Abide by Crawler delay in Robots.txt

# Web Crawler Parallelisation Policy

- An efficient crawler needs to access many web servers at once

- Run multiple processes in parallel

- Could find the same URL on multiple pages

- Strategies:

    - Dynamically assign pages to crawler processes

    - Static mapping e.g. based on a hash function

# **Web Crawler** Crawlability

- Broken links

- Denied access

- Outdated URLs

- URL errors

- Blocked

- No out links

- Slow load speed

- Duplicated pages

- Flash content

- JavaScript (Googlebot executed from 2014)

- HTML frames (outdated and thus poorly indexed)

# **Search Index** Ranking

- Query results can be ranked using many features:

    - How many times does the page contain the keywords

    - Do keywords appear in the title or URL

    - Does it contain synonyms for your keywords

    - Is it from a quality source

    - What is its importance

    - How often a page is updated

    - Freshness of information

    - Page load time

# Indexing Data Issues

- Distributed data

- Changes over time

- Large volume

- Data

  – Unstructured data - gifs, pdf, etc

  – Redundant data - 30% pages are near duplicates

  – Quality of data - False, poorly written, invalid, misspelt

  – Heterogeneous data - media, formats, languages, alphabets

# User Search Problems

- Users do not understand how to provide a sequence of words for searches

- Users may get unexpected answers because they are not aware of the input requirement of the search engine.

  - For example, some search engines are case sensitive.

- Users have problems understanding Boolean logic

- Around 85% of users only look at the first page of the result, so relevant answers might be skipped.

# User Search Problems: Ordering of Terms

# User Search Problems: Ordering of Terms

# Boosting web traffic

- Search engines direct traffic to websites

- 85% of people don't look at the second page

- People try to optimise their site so that it ranks highly on Search Engines

- Could be fundamental to a website's business model

# Search Engine Optimisation

- Whole industry exists trying to boost search ranking to ensure pages are indexed by search engines

- Legitimate SEO (White Hat)
  - Good Design
  - Valid metadata, alt tags on images

- Illegitimate SEO (Black Hat)
  - Often gaming search ranking algorithms
  - Deception
  - Leads to arms race between SEO and search engines

# Combatting SEO

- Most search engines have rules against:
  - Invisible text
  - Meta tag abuse
  - Heavy repetition
  - "domain spam"
    - Overtly submission of "mirror" sites in an attempt to dominate the listings for particular terms

# Google and other SEs are a Business

- Search Engines record tracking information
  - Google saves every voice search
  - IP addresses
  - Location
  - Saves your searches
- Google's revenue is from adverts
  - Improve their revenue with targeted advertising
- Google has a large research department
  - Improve their technology

# History of the Web Browser

# World Wide Web

The WorldWideWeb (W3) is a wide-area hypermedia information retrieval initiative aiming to give universal access to a large universe of documents.

Everything there is online about W3 is linked directly or indirectly to this document, including an executive summary of the project, Mailing lists , Policy , November's W3 news , Frequently Asked Questions .

What's out there?
>    Pointers to the world's online information, subjects , W3 servers, etc.

Help
>    on the browser you are using

Software Products
>    A list of W3 project components and their current state. (e.g. Line Mode ,X11 Viola , NeXTStep , Servers , Tools , Mail robot , Library )

Technical
>    Details of protocols, formats, program internals etc

Bibliography
>    Paper documentation on W3 and references.

People
>    A list of some people involved in the project.

History
>    A summary of the history of the project.

How can I help ?
>    If you would like to support the web..

Getting code
>    Getting the code by anonymous FTP , etc.

# Before Search Engines

- Exploit hyperlink structure

  – Personal homepages with links

  – Directories

- Word of mouth

  – Email, forums, Usenet

# Timeline 90's – Hand vs Crawler

June 1993
1st Web Robot
Wandex
**Meausres size**

Sept1993
1st Web SE
**By Hand**

Oct/Nov 1993
2nd Web SE
**By Hand**

Dec 1993
3rd Web SE
**Crawler**

Jan 1994
Altavista
**Crawler**

April 1994
Yahoo!
**Web Directory
By Hand**

Jul7 1994
Lycos
**Crawler**

1995
LookSmart
**Web Directory**

Marh 1998
Larry Page
BackRub
SE

May 1996
HotBot
SE

April 1997
Ask Jeeves
SE

1998
Google Search
SE

Jul/Sept 1998
MSN Search

- Services to search the web started in the early 90s
- The indexes were either created and edited by hand
- Or, by crawlers

1993          1994          1995          1996          1997          1998

# Early 90s Architectures

# Early 90s Architectures



Query Engine

Interface

Search index

Web Crawler Or Person

# Web Crawlers

- **Wandex** 1993 – size of the web

- WebCrawler Dec 1993 - indexer

- *WebCrawler* 1994 - indexed entire web page


WebCrawler

WEBCRAWLER™

To search the WebCrawler database, type in your search keywords here. Type as many relevant keywords as possible; it will help to uniquely identify what you're looking for.

Pets

Search          ☒ AND words together

Number of results to return: 25 ▼

# Yahoo! 1994

# Welcome to PizzaNet!

PizzaNet is Pizza Hut's Electronic Storefront and is brought to you by Pizza Hut®and The Santa Cruz Operation® You may click on the Pizza Hut logo on any page to submit comments regarding PizzaNet to webmaster@Pizzahut.COM

---

If you would like to order a pizza to be delivered, please provide the following information:

**Name:**

**Street Address:**

**Voice Phone ###-###-####**

(where we can reach you)

Continue

# Website - Telegraph 1994

# Website – New York Times 1995

# Website – BBC News 1997

# Number of Websites June 1994 - Dec 1995

The original Amazon website (August 1995)
Source: Restored by Taran Van Hemert

# Search Engines Issues Mid '90s

- Scale
  - How to rank pages to give the best results

- Spamming
  - Web searches favoured web pages with high keyword density

- Keywords
  - Spelling mistakes

# Timeline 90's

June 1993
1st Web Robot
Wandex
Meausres size

Sept1993
1st Web SE
List

Oct/Nov 1993
2nd Web SE
List

Dec 1993
3rd Web SE
Crawler
Indexer
Seraching

Jan 1994
Altavista

April 1994
Yahoo!
Web Directory

Jul7 1994
Lycos
SE

1995
LookSmart
Web Directory

March 1996
Larry Page
BackRub
SE

May 1996
HotBot
SE

April 1997
Ask Jeeves
SE

1998
Google Search
SE

Jul/Sept 1998
MSN Search

| 1993 | 1994 | 1995 | 1996 | 1997 | 1998 |

# BackRub

- Utilised PageRank
  - More advanced than previous indexing
- Used back links

# PageRank

- Assumes important websites are likely be linked to

- Similar to a voting:

  – where in links are votes

  – the quality of a vote is determined by the number of votes (in links)

  – your vote is worth more if you have a higher page rank

- Factors Considered:

  – Number of in links

  – Quality of in links

  – Web page's context

- Favoured high keyword density

# Keyword Stuffing

Our company provides software development and consultation services.  We are the best software development and consultation services. Our company staff provides software development and consultation services. We offer the best software development and consultation services. Our software development and consultation services is the right solution for your business. Satisfaction with our software development and consultation services is guaranteed. We build your business solution with out software development and consultation services.  Contact us and you will receive quality software development and consultation services.

# Keyword Stuffing

Our company provides **software development and consultation services**. We are the best **software development and consultation services**. Our company staff provides **software development and consultation services**. We offer the best **software development and consultation services**. Our **software development and consultation services** is the right solution for your business. Satisfaction with **our software development and consultation services** is guaranteed. We build your business solution with **out software development and consultation services**. Contact us and you will receive quality **software development and consultation services**.

# The Original Google Storage

- Larry Page needed a large amount of diskspace to test PageRank

- 10 4GB hard drives

# Google 1997-8

# The Evolution of Search Engine Features

- Search Engine algorithms are constantly changing

- They strive to improve their:

  - Efficiency

  - How to order the results

    - Freshness

    - Relevancy

  - How to improve user interaction

  - Additional features

# Google Adwords 2000 Funding

# Relevancy

- During 9/11 it was apparent that search engines did not cater of time-sensitive queries

- This lead to Google developing Google News

- The freshness of search engine's index became more important

# Google 2012

- In 1999, it took Google one month to crawl and build an index of about 50 million pages.*

- In 2012, the same task was accomplished in less than one minute.*

- 16% to 20% of queries that get asked every day have never been asked before.*

*Mitchell, Jon. "How Google Search Really Works." Readwrite. February 29, 2012.

# Search Verticals

- In a bid for content and market share, Search Engines develop other search verticals:
  - Books
  - Travel
  - Finance
  - Shopping
  - Scholar

# Mobile Search

- By 2015 mobile accounted for more than half of all searches

- Searches need to be context aware (location)

- Limited screen and bandwidth means relevancy is critical

- Boost mobile-friendly pages in ranking algorithm

# Next: Web Standards