# How Google Works:

— Coding sloth

Search Engine

① Crawling            ② Indexing                    ③ serving
(grab the web
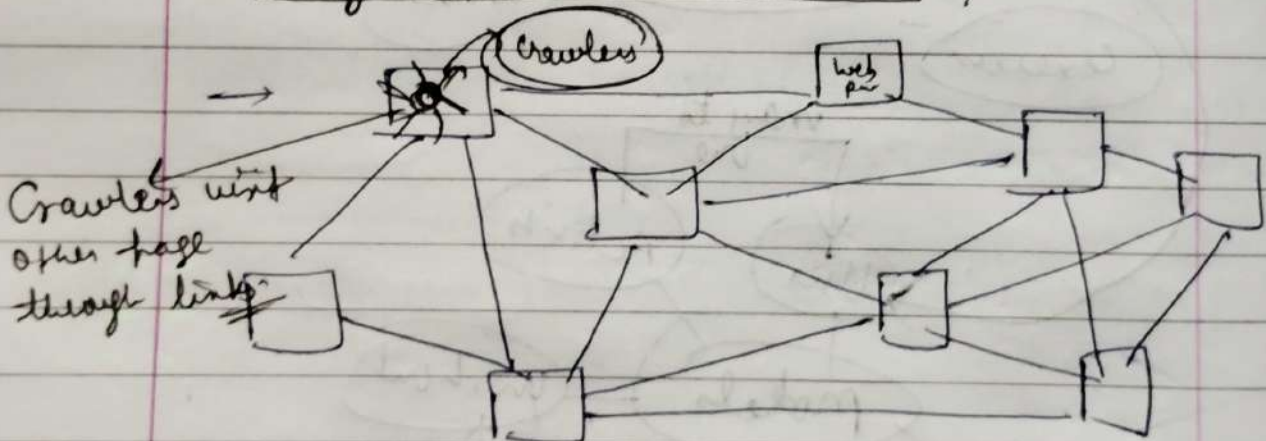page)

① Crawling.

→ through crawling google knows every website you visit.

→ automated programs: Crawlers / Spiders

⟹ why is internet called web?

→

Crawlers
web
pa

Crawlers visit
other page
through links

each web page linked to each other
→ looks like spider web

→ URL discovery

(Sitemap)

↓

file which keep information
about pages

→ Google has bunch of crawler

Google bot
Desktop

Google bot
mobile

(Google bot) ──┌─────────────────────────────────┐
              │ → (1) which page to crawl       │
              │ → (11) How often it should be crawled │
              │ → (11) How many pages to crawl  │
              └─────────────────────────────────┘

Algorithms

(robots.txt)
↳ this file tells which webpage bot
should crawl.

Avuoy

Page No.:
Date:

Page No.:
Date:

youvA

## Our Bot

→ it visits the url (our case: wikipedia page of Google)

↪ then it will grab all links in webpage

→ then all the grabbed links will be put in list to be visited later

→ we will then move on to new url and repeat the process

time sleep
Thread.sleep } make sure we donot overload the Server

↓

if it detects that it is bot

↓

Banned!

Controller → service

## ① J soup

↳ API for extracting & manipulating data
↳ deal with real world HTML

→ Fetching url : Jsoup.connect ("https://example.com").get()

→ Selecting elements :
   Elements links = document.select ("a[href]");

→ Extracting element attributes
   for (Element link : links) {
      String href = link.attr("href");
      String text = link.text();
   }

→ Manipulating
   Element element = document.select ("div.classname").first();

   element.text ("New text");

→ CRAWL LIMIT :- limit of pages to visit

Set                         Queues
↳ track of              ↳ mange URL
  URL visited              to be crawled
↳ avoid                      ⇓
  recrawling              [poll] remove

→ Thread·sleep (1000 + random·nextInt (2000))
        ↳ delay b/w 1 & 3 sec
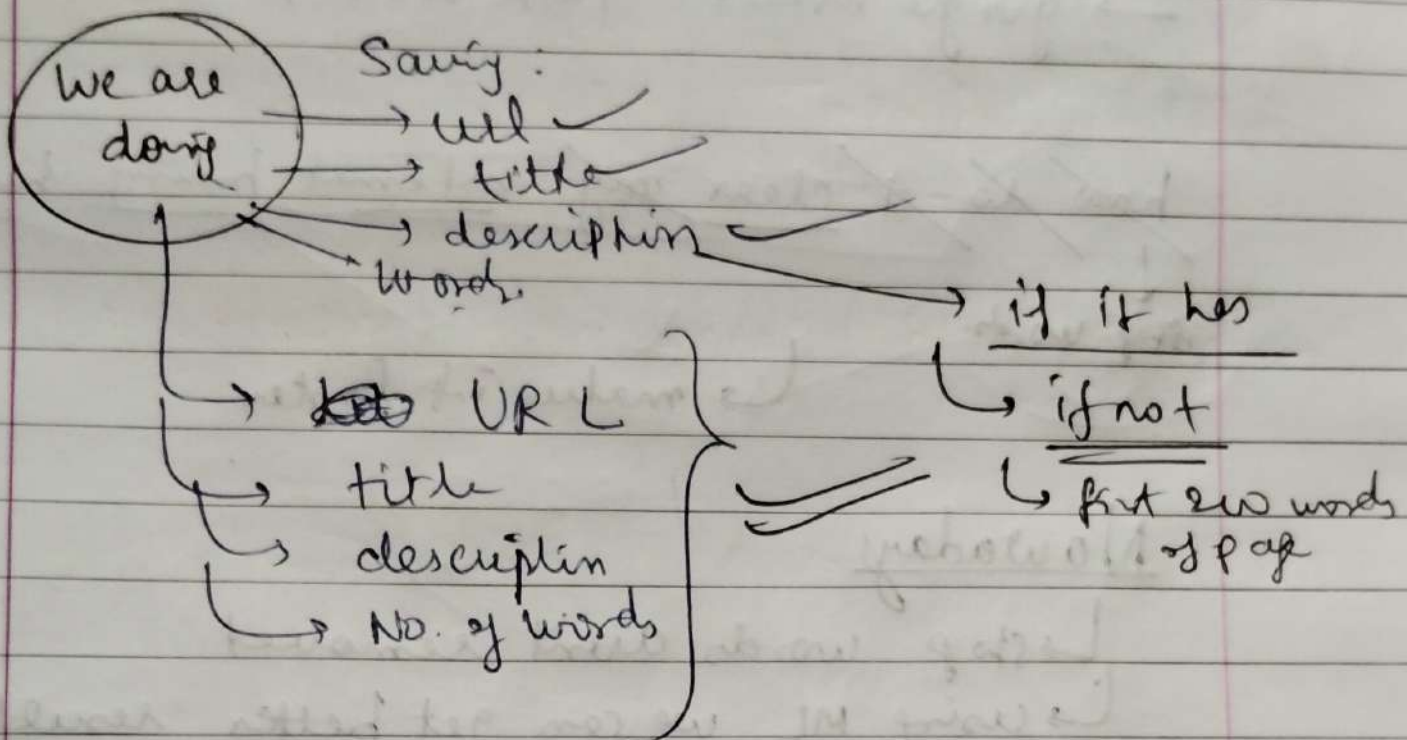
→ fetch Document

→ format Url { null
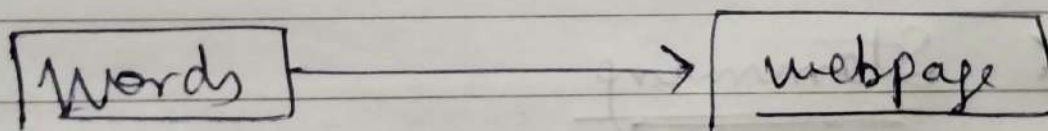               invalid

DOES WHAT IT NEEDS TO DO!
WHO KNOWS WHETHER
              CORRECT ?

# (2) Indexing

→ google saved every info from website.

We are doing → Saving:
→ url
→ title
→ description
→ words.

→ URL
→ title
→ description
→ No. of words

→ if it has
→ if not
→ first 200 words of page

## ★ Inverting Indexing

| words | → | webpage |

| Java Program q |
→ DB containing 'Java Prog'  words
→ return the websites having word 'Java Prog'

## ⭐ Stop words

the, is, at, which, on

→ google remove these words

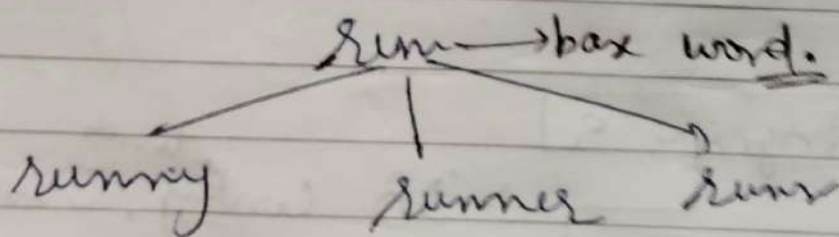how do I clear my internet browsing history?
↓
stop words.
        ↳ makes it faster

### Nowadays
↳ stop words aren't removed
↳ using ML we can get better results
   through stop words

## ⭐ Stemming

Run → base word.
        ╱    |    ╲
   running  runner  runs

Searching this will consider base word

## Indexing Service

**model**

1. Collect title
   - document title

2. description
   - if it has ~
   - it not : first 200 words

3. words — count

   extract

→ Pattern  } [List]
→ matches

**url**
**title**
**description**
**words**

(getter/setter)

## Web Crawler Controller

[map]
→ indexing service

③ Serving



webpy (A)

Page Rank
Algorithm

0

(B)

3
webpage

dsk

webpage
(D)

0 2 (C)
webpy

→ Topological Sorting

Ranking:  B > C = D > A

$$PR(P_i) = \frac{1-d}{N} + d \cdot \sum_{P_j \in n(P_i)} \frac{PR(P_j)}{|out(P_j)|}$$

WTF?