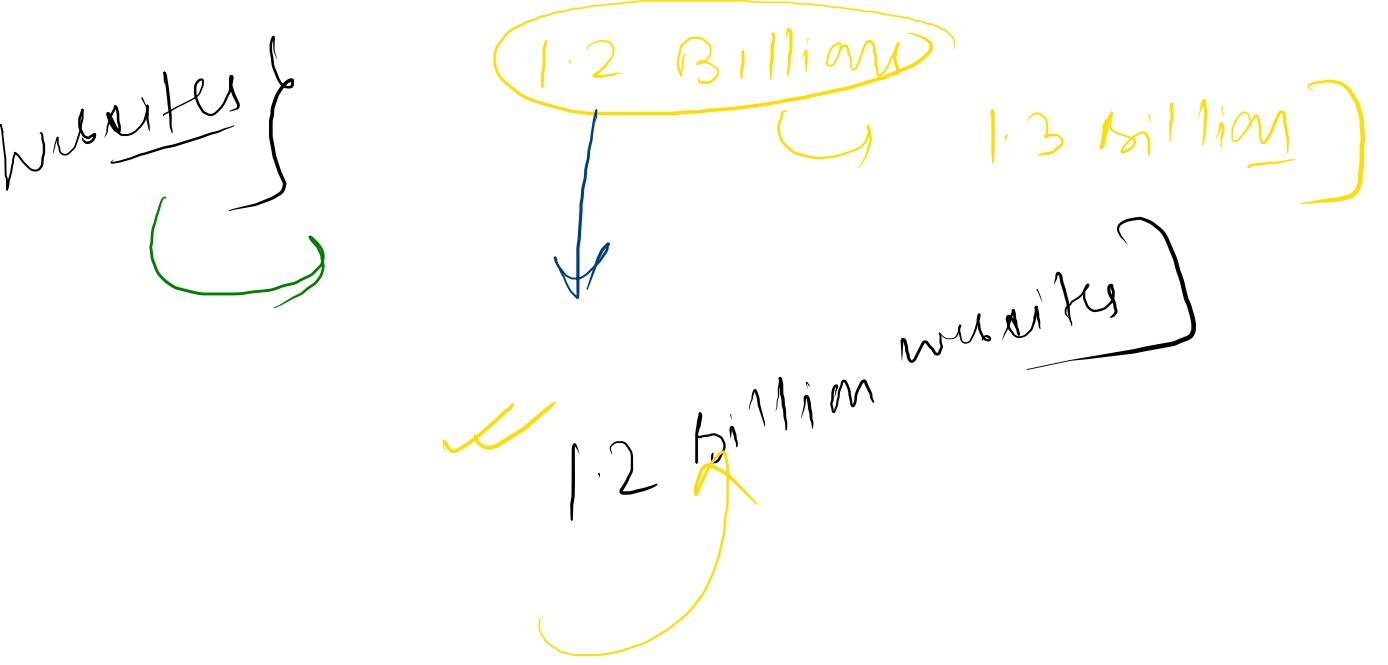
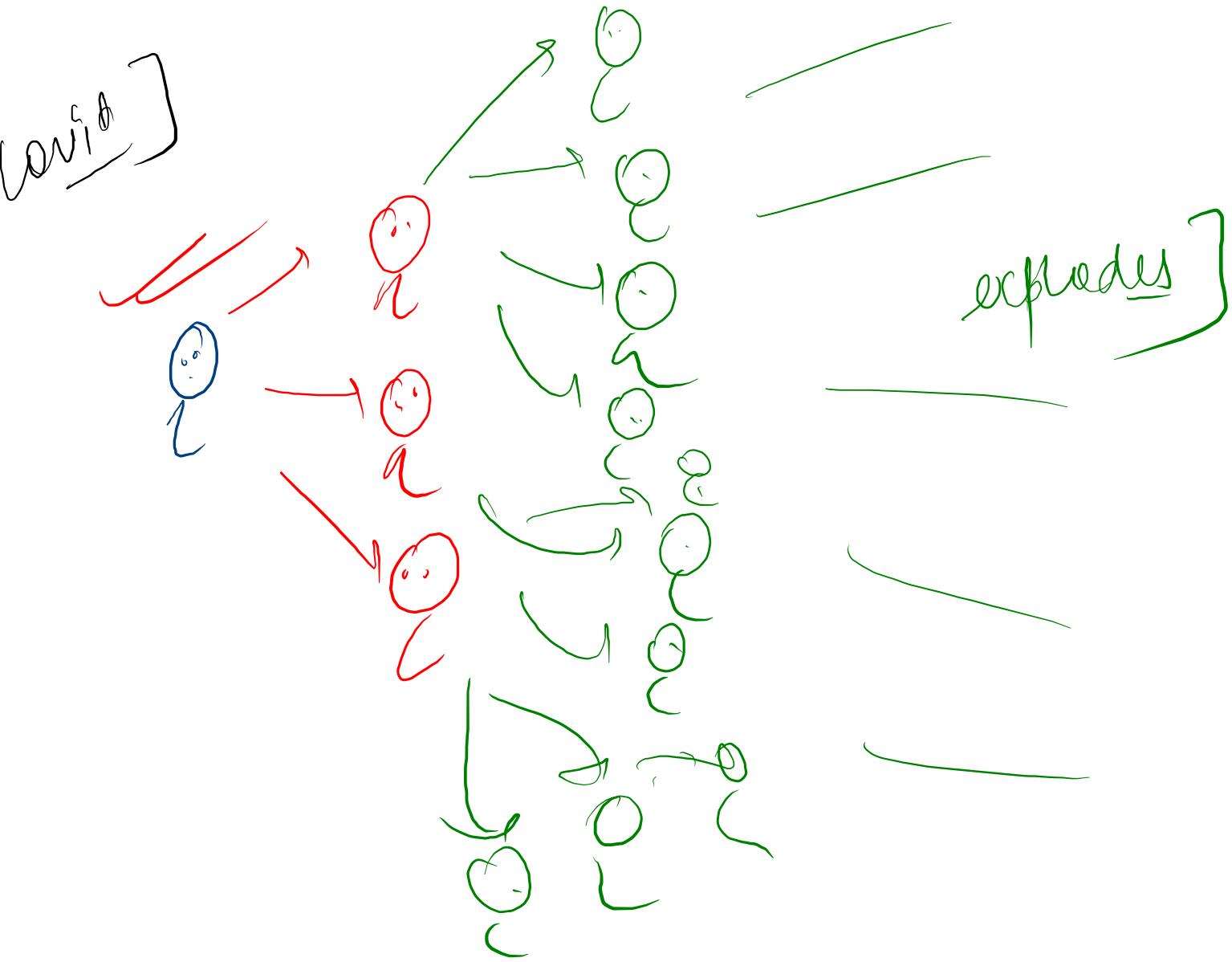
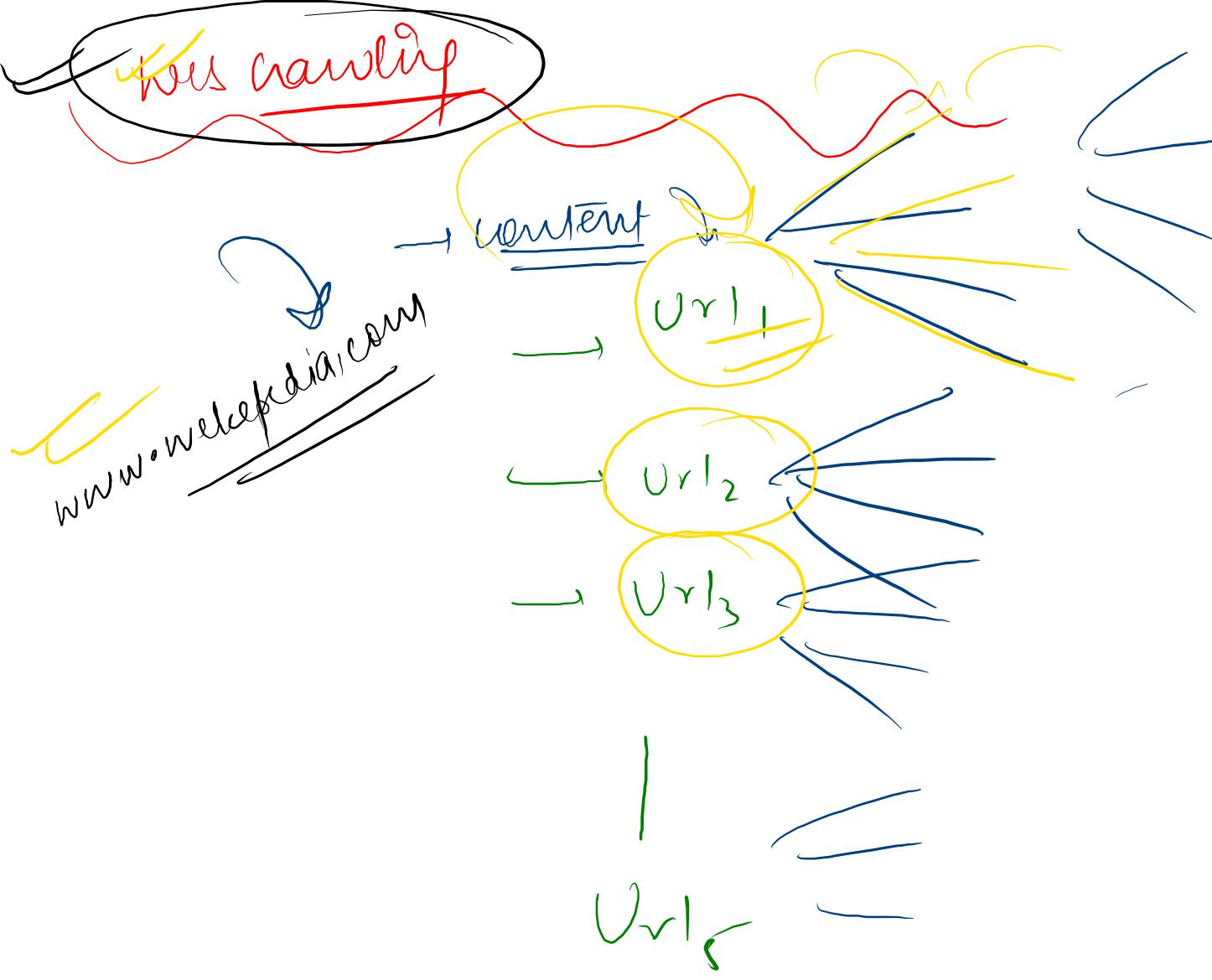


Web crawler?



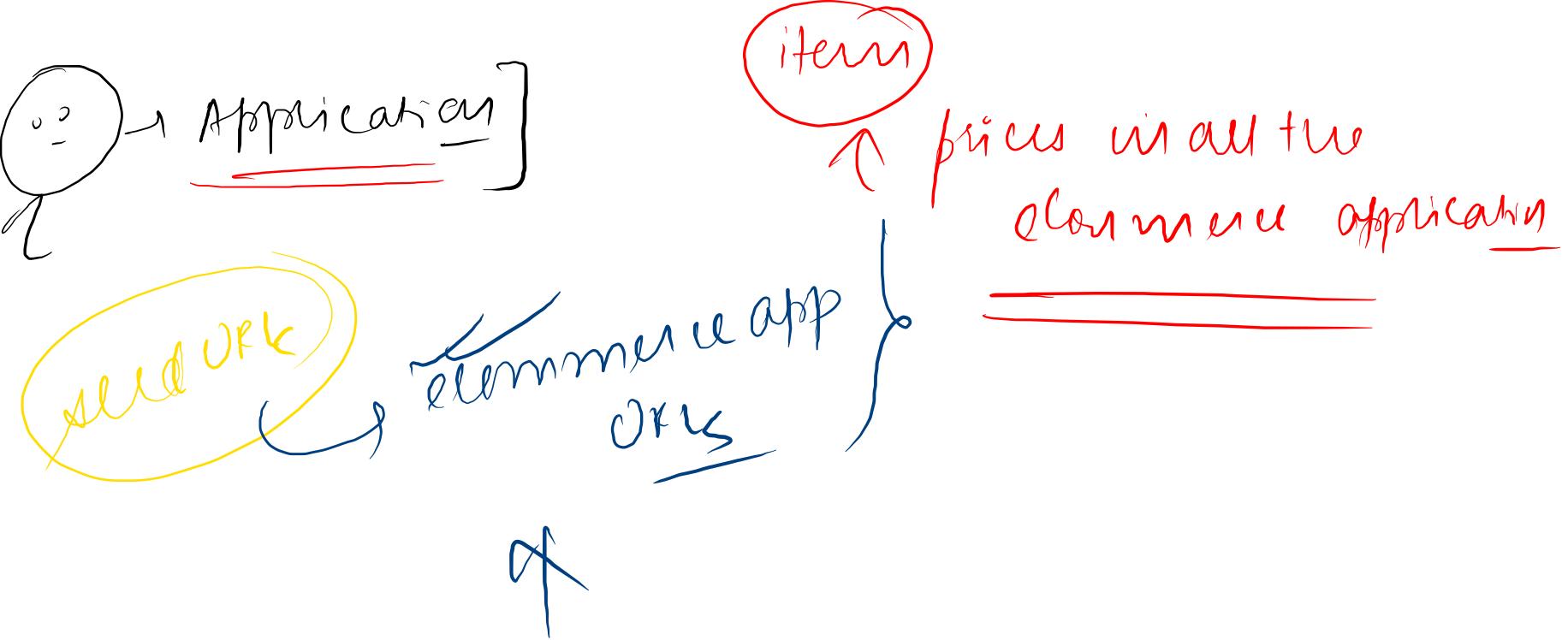






some basic website links

see URLs



## Web Crawling

- ① Search engine → Google → Bing
- ② Copyright violation detection/  
Identify the copying of reports  
plagiarism
- ③ Data source → Data fetching
- ④ Keyword based search

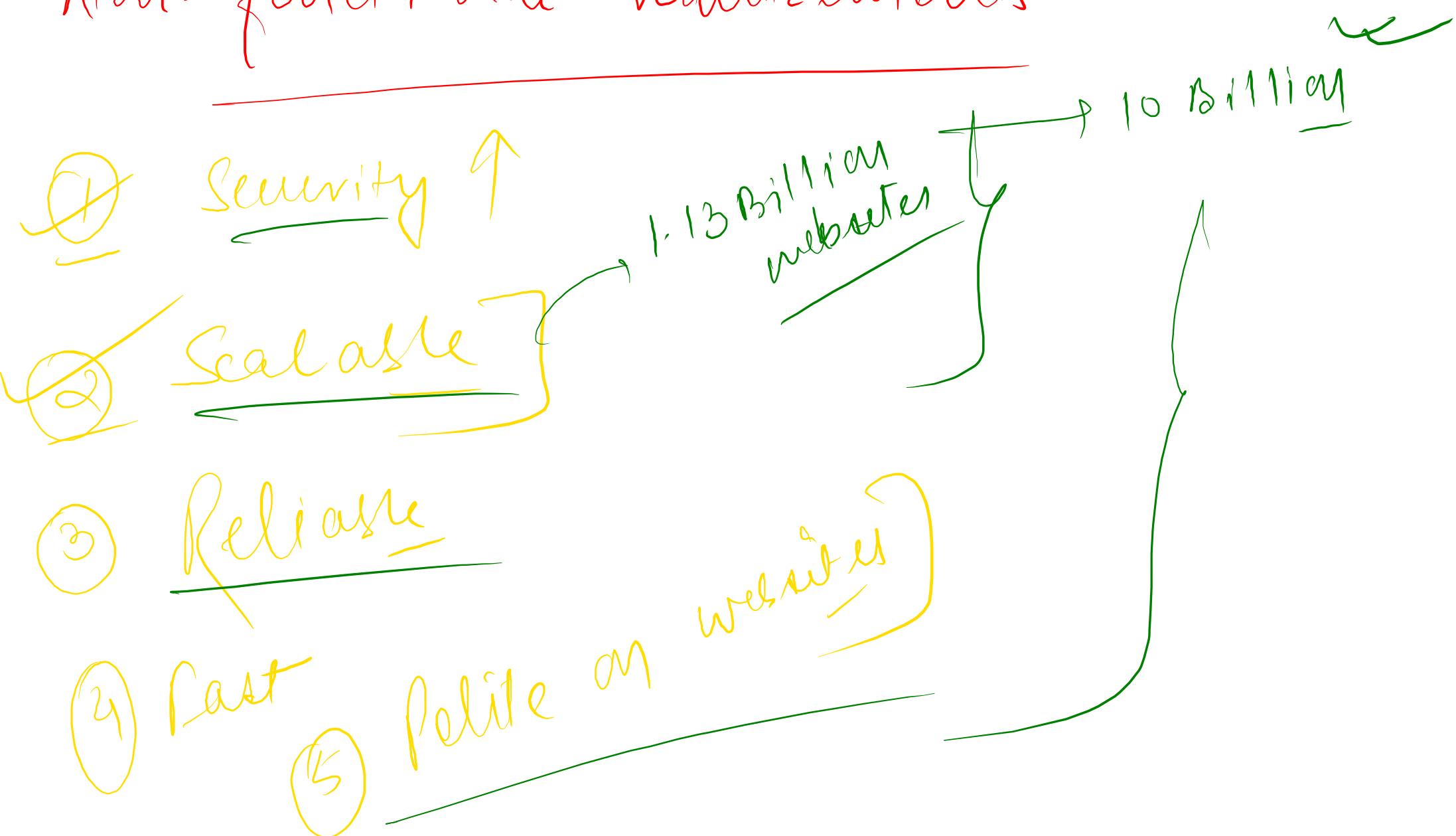
:-> Build this system

## Functional Requirements

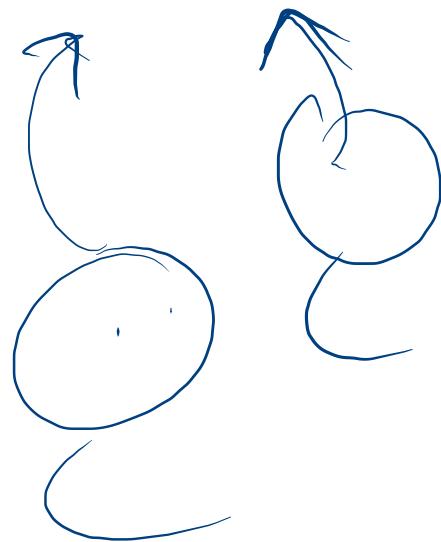
① Scan and store data for all the websites

② Don't store duplicate records

# Non-functional requirements



user → Website ] public content ]



Estimation  
↳ estimate the scale of the problem

X

of pages

Total no of website = 1.2 Billion

60% of website are active

$$\therefore \frac{60}{100} \times 1.2 \text{ billion} = 720 \text{ million}$$

Facebook  
Wikipedia



Average  $\rightarrow$  No of web pages / website = 100

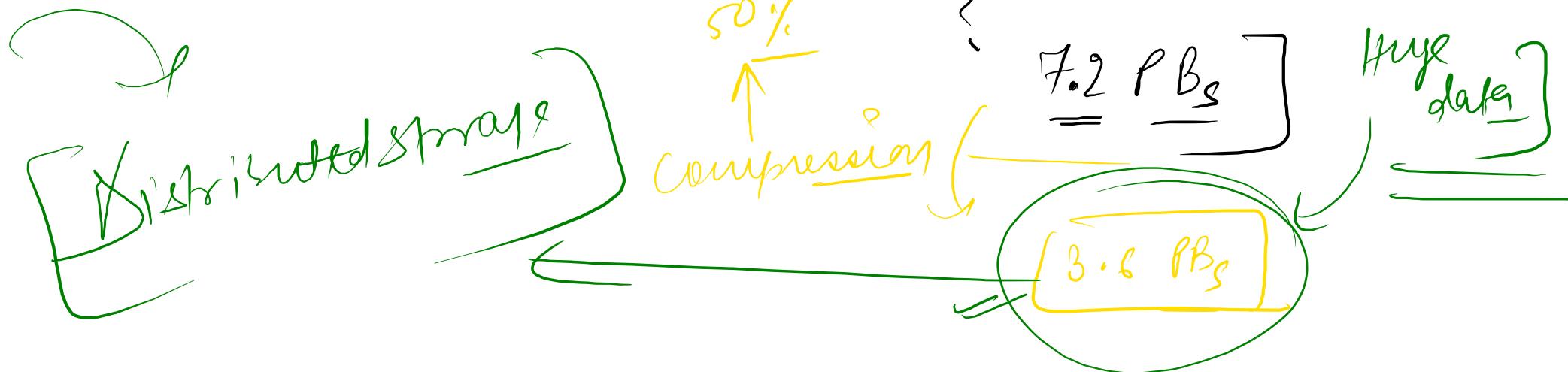
$$\therefore \text{Total no of webpages/URLs} = 720 \text{ million} \times 100 \\ = 72000 \text{ million} \\ = 72 \text{ billion}$$



Avg size of 1 webpage = 100 kB

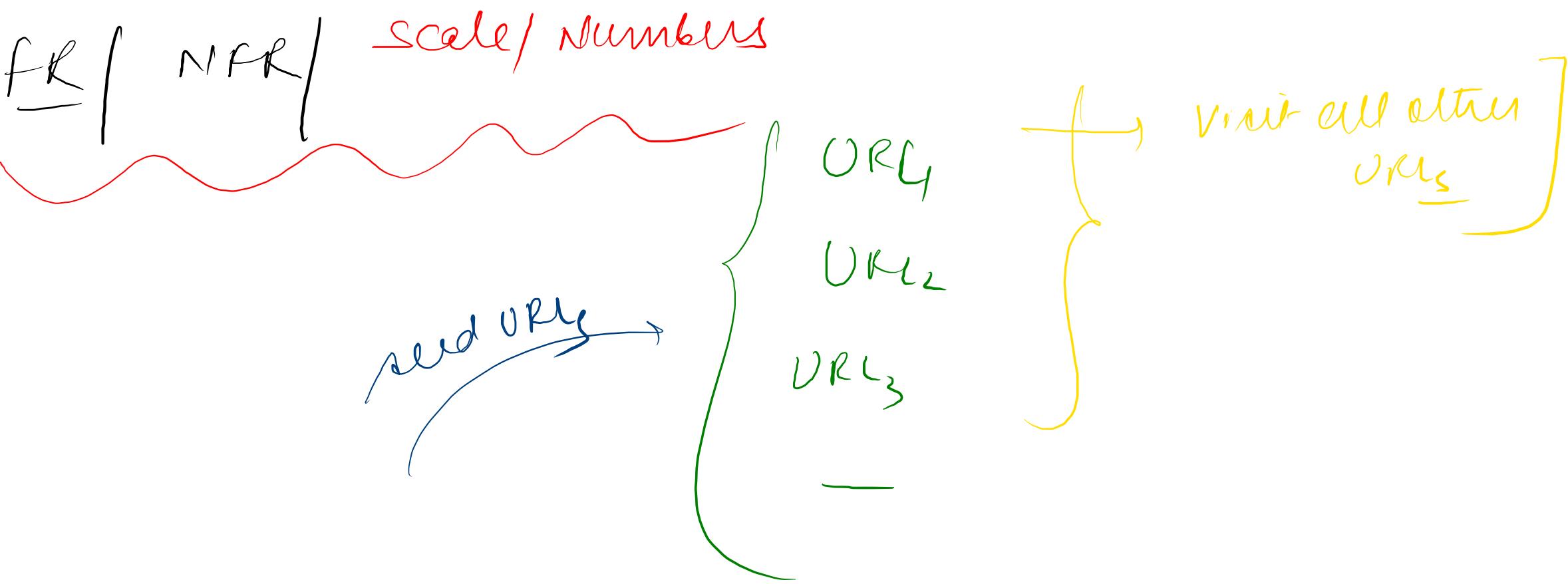
No of pages = 72 billion

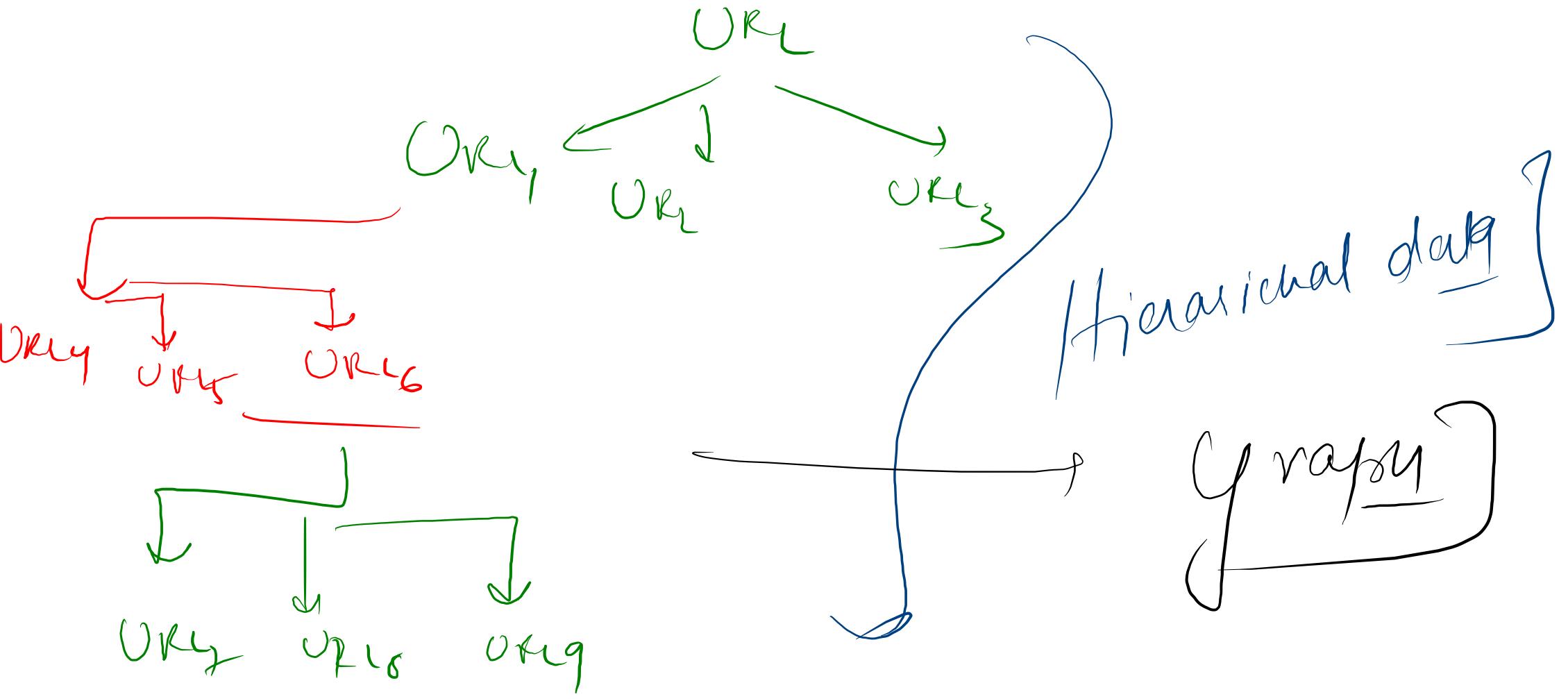
total size of webpage to be stored  $\approx 72 \text{ billion} \times 100 \text{ kB}$



raw  
multiple  
line

→ updating the data  
spread

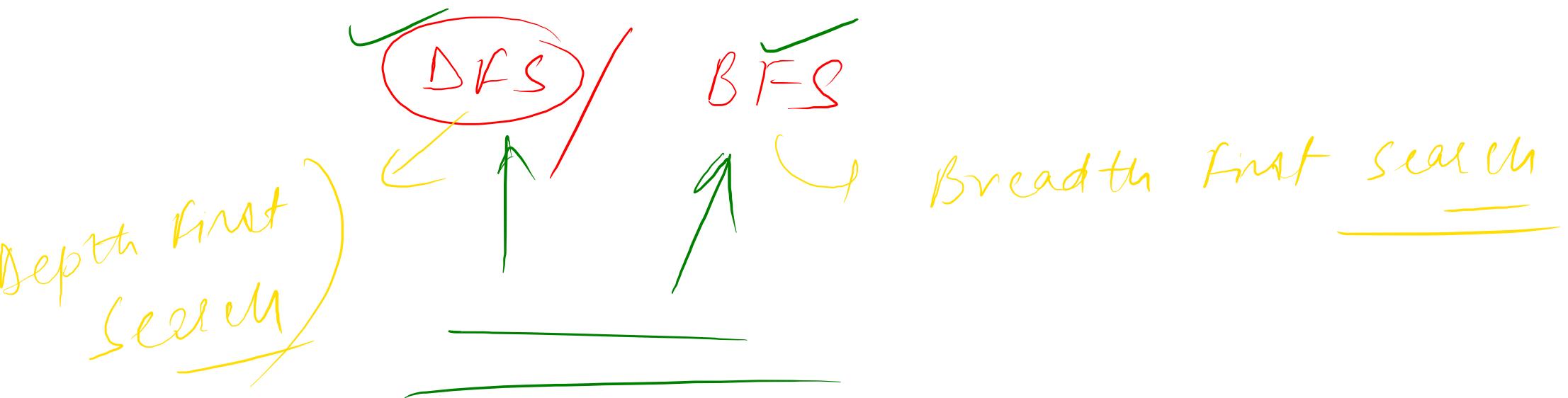




Complex graph

A Undirected graph

Graph traversal



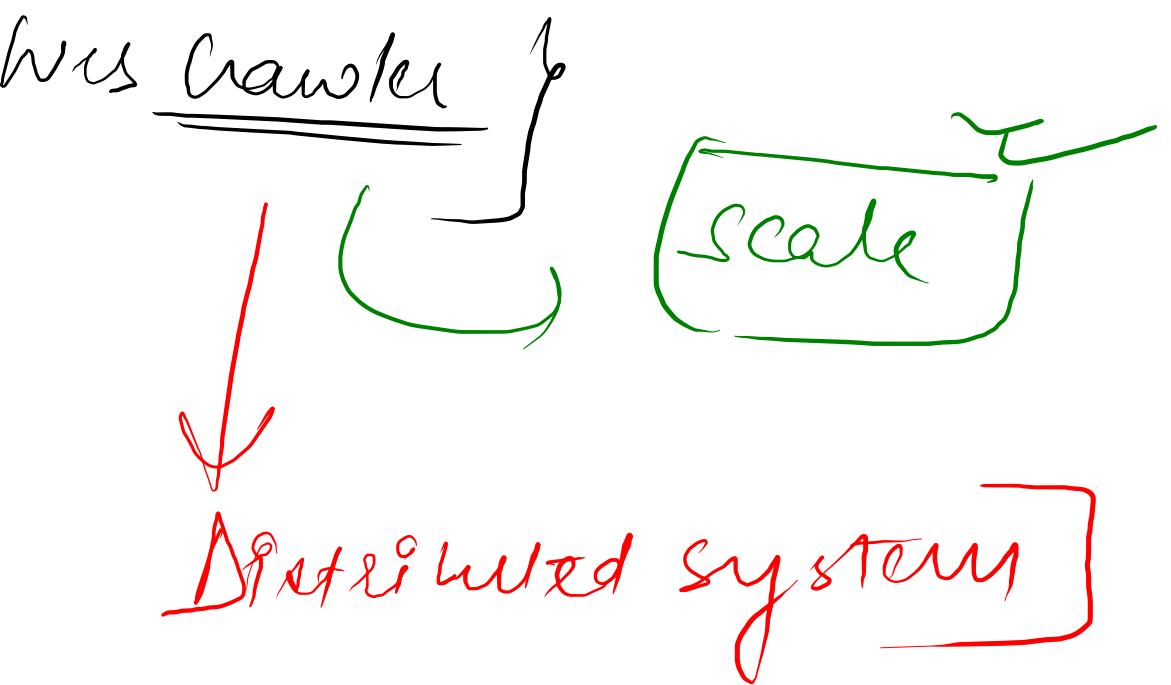
BFS

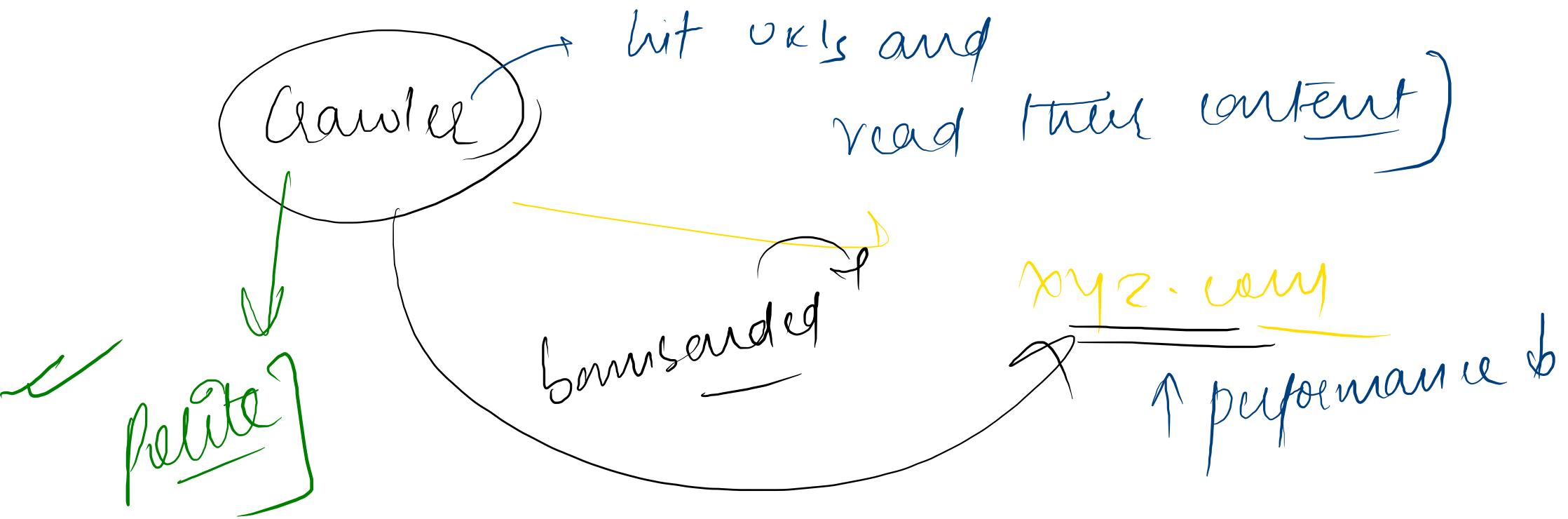
?

(DFS) ?

~~Both q time return all possibilities  
eventually~~

The diagram illustrates the flow of data from memory to storage. A large bracket labeled "Cache" covers both "RAM" and "LTB". An arrow points from "Cache" to "String data in-memory". Another arrow points from "RAM" to "LTB".





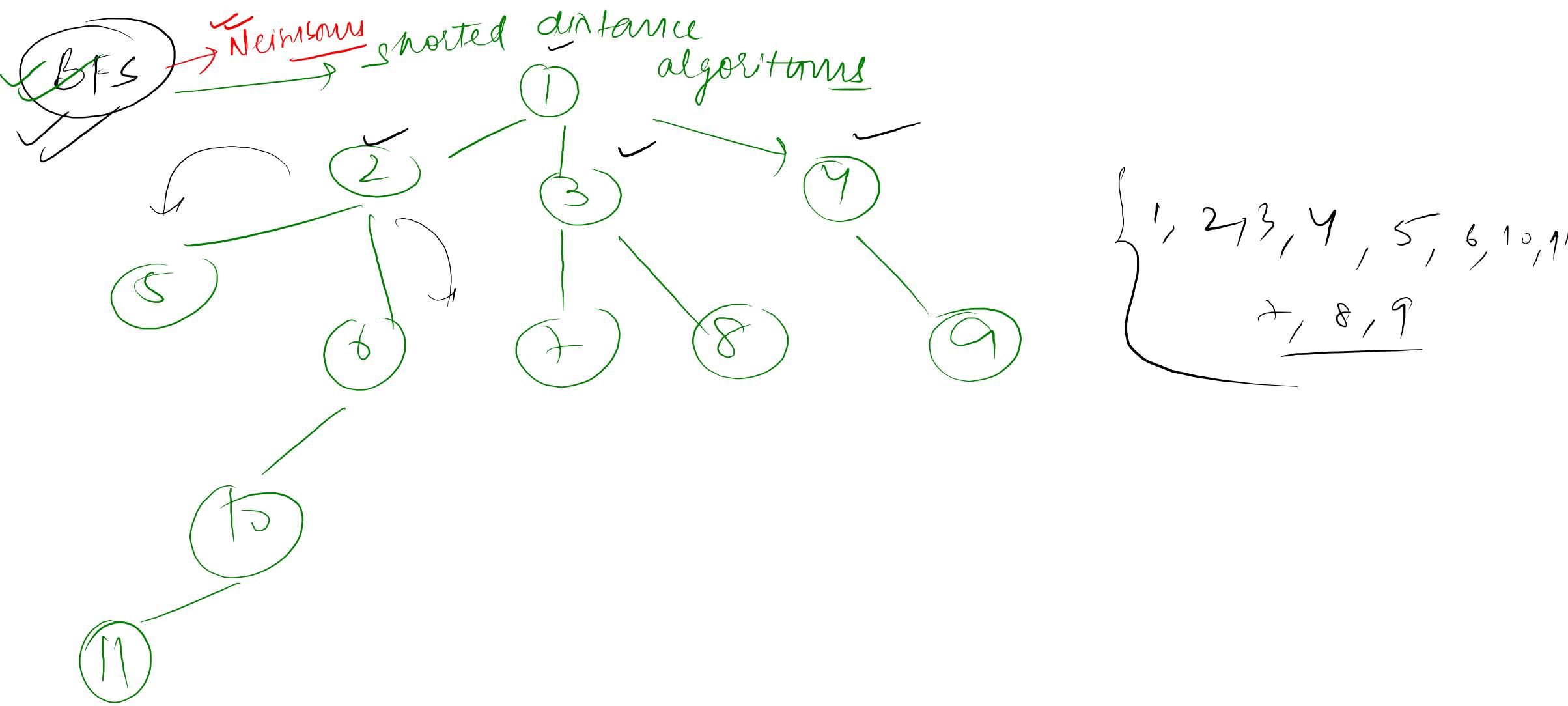
what's

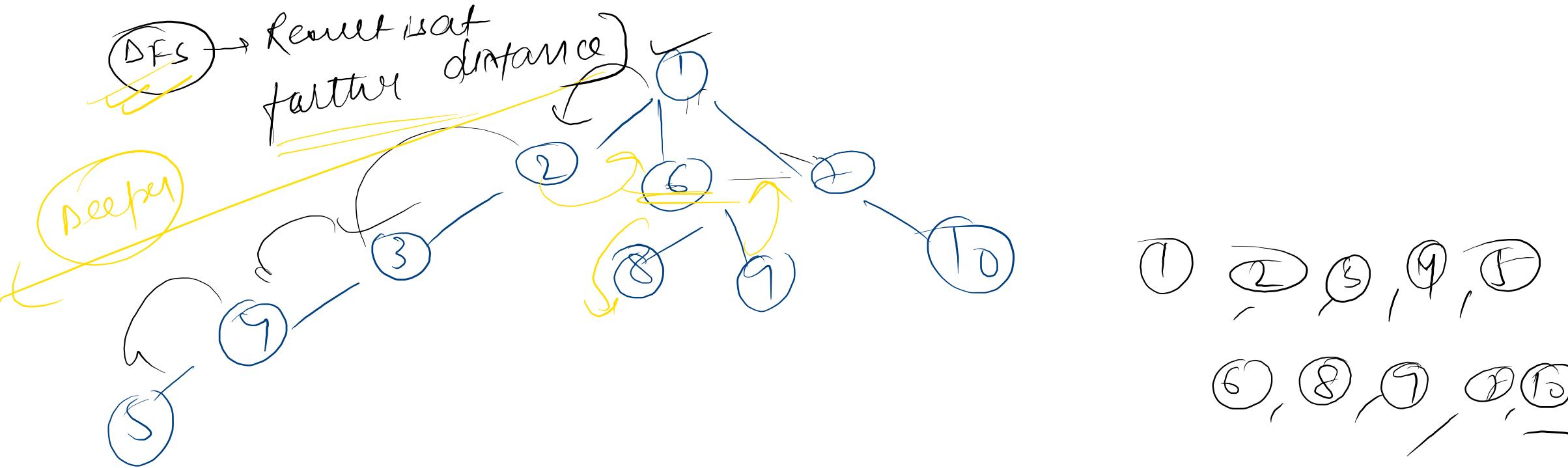


name should

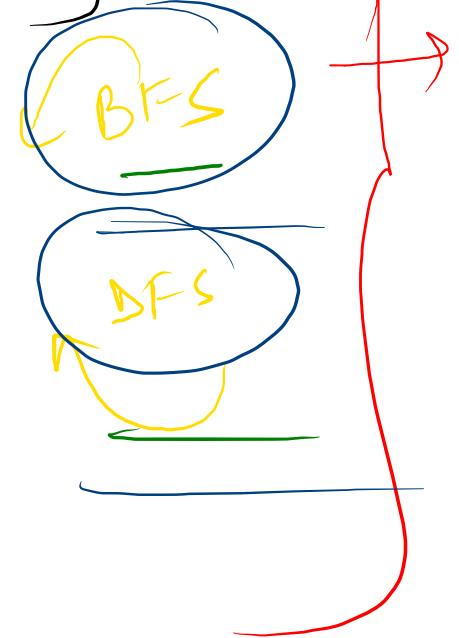
A green bracket on the left encloses the phrase "name should". A red bracket on the right encloses the word "index". A green curved arrow points from the green bracket to the red bracket.

what/why } we are looking ]



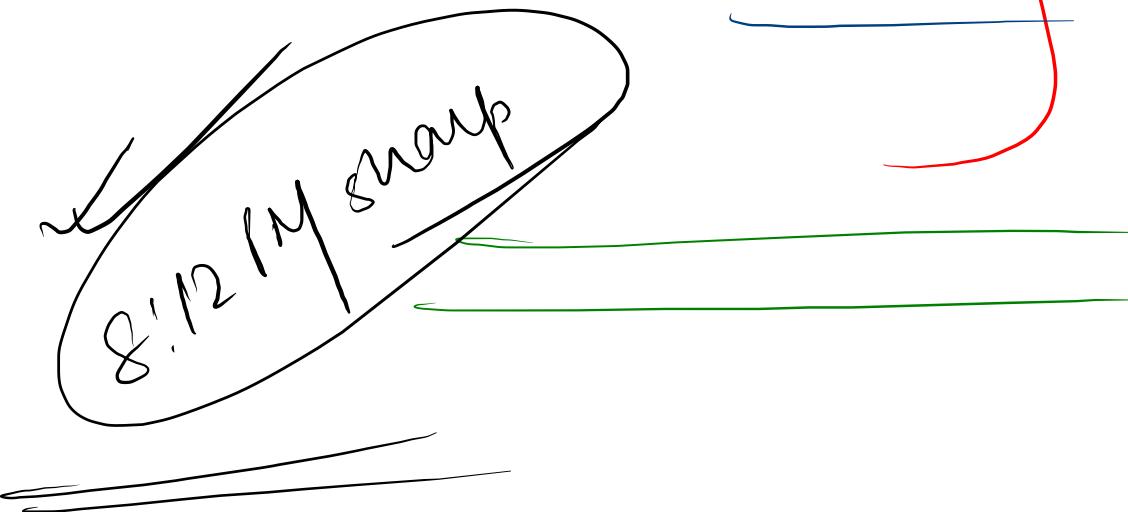


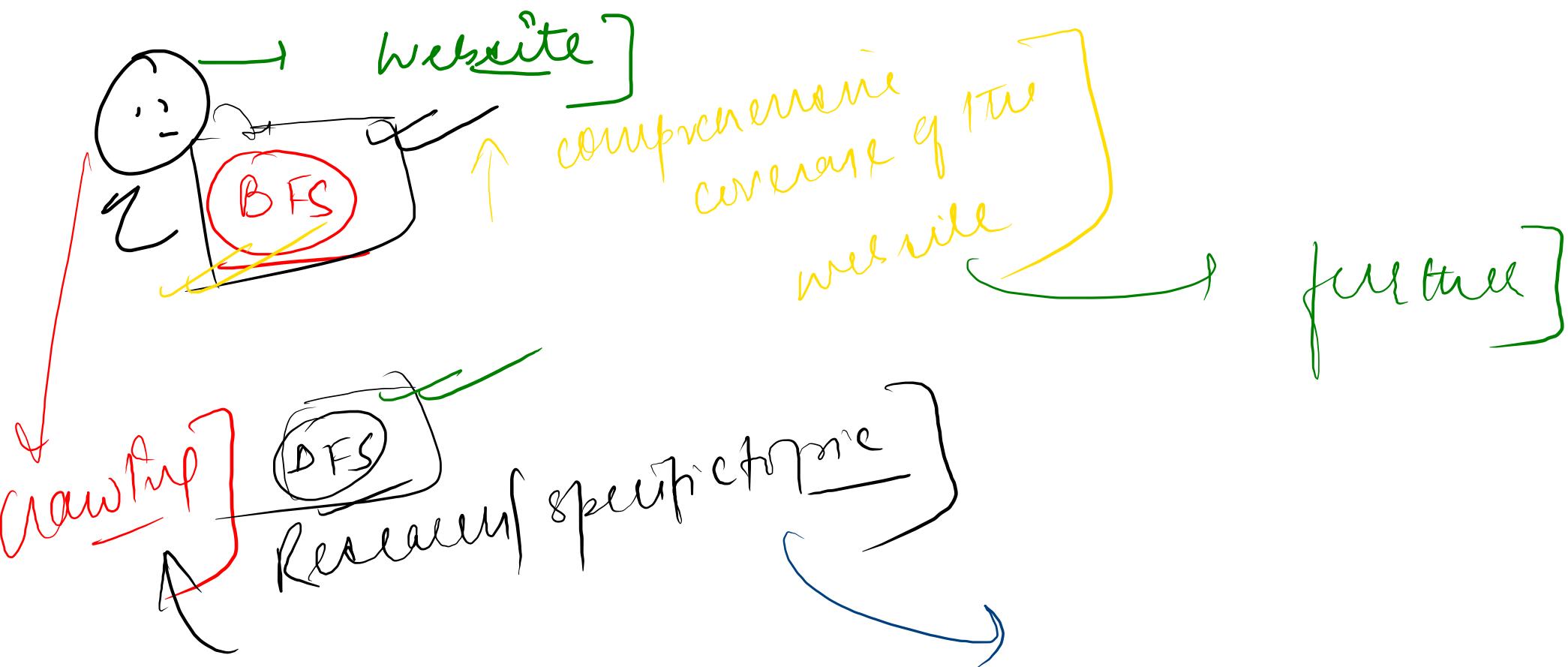
BFS and DFS



Prefer? / why?  
when

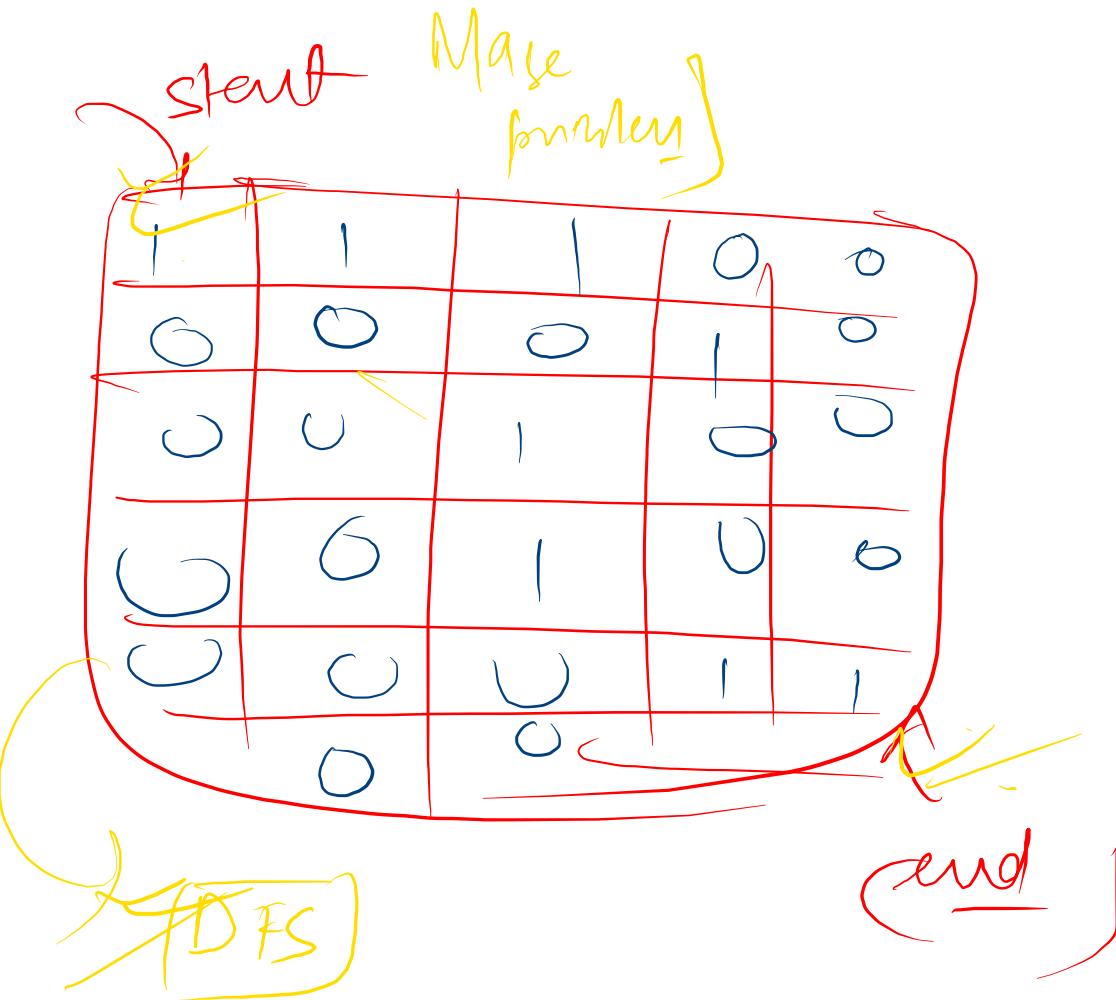
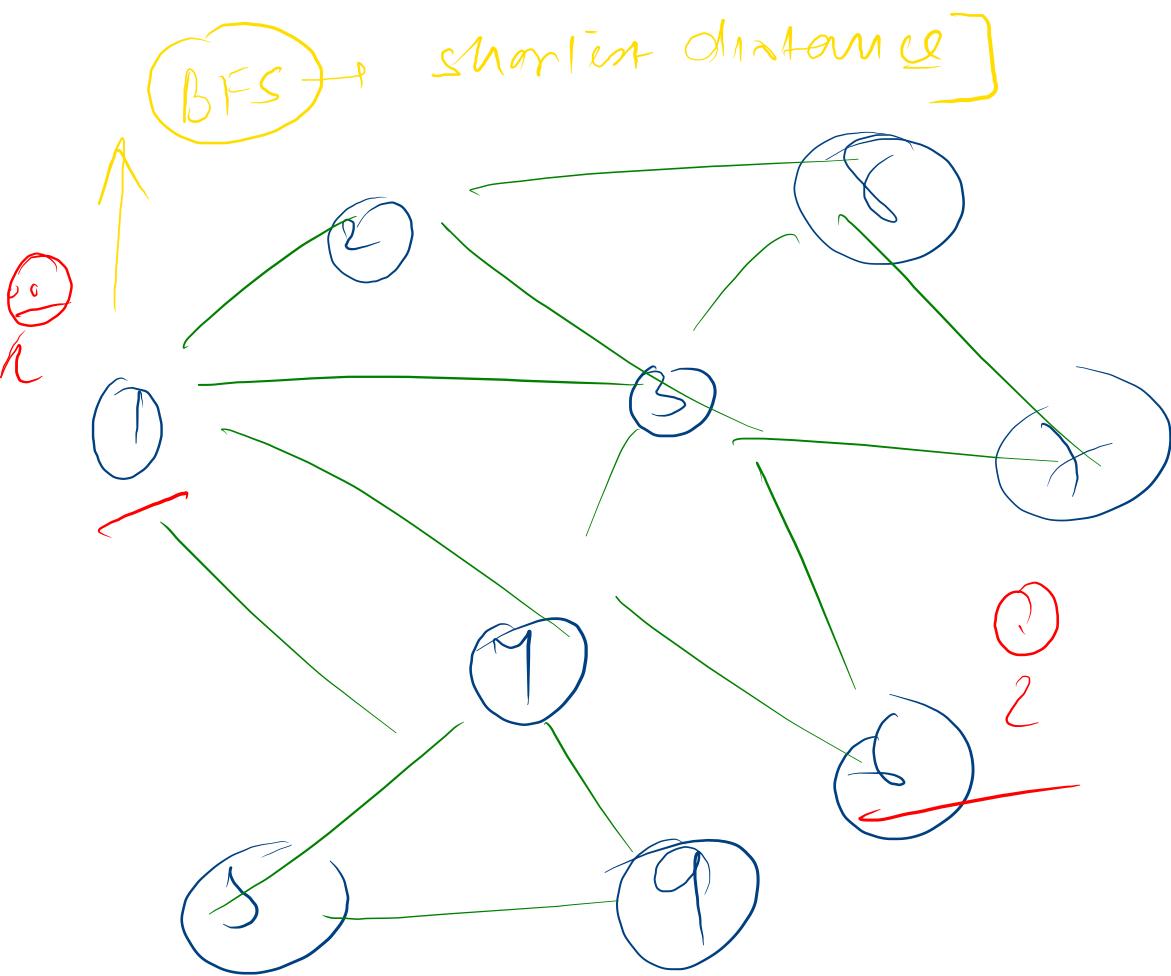
found break





~~BFS~~

scanning / crawling all the  
websites



Google Search

indexed results

→

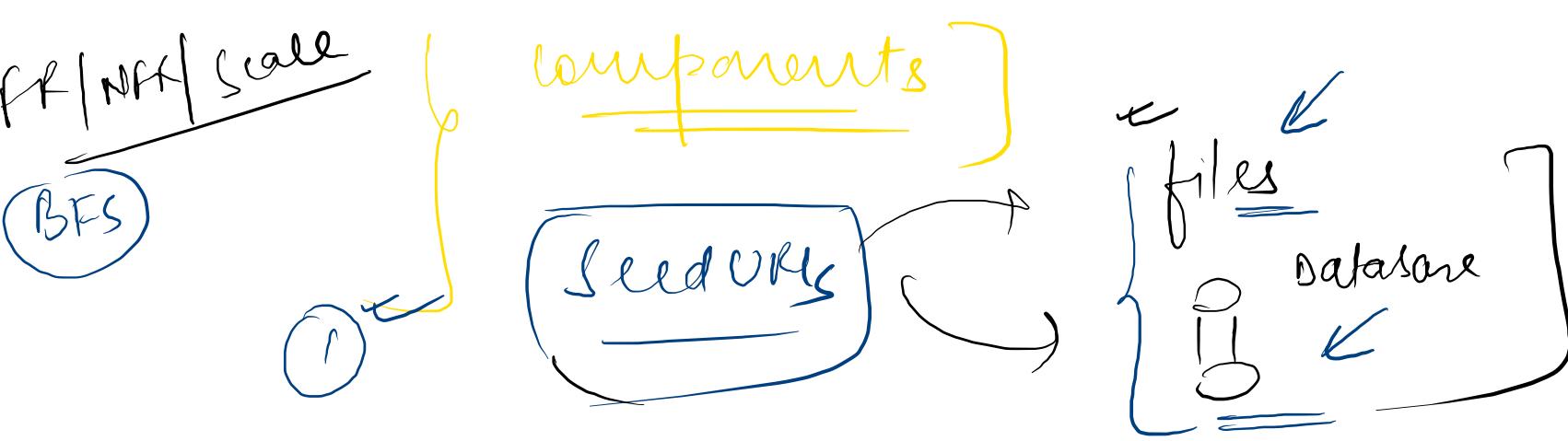
stored data

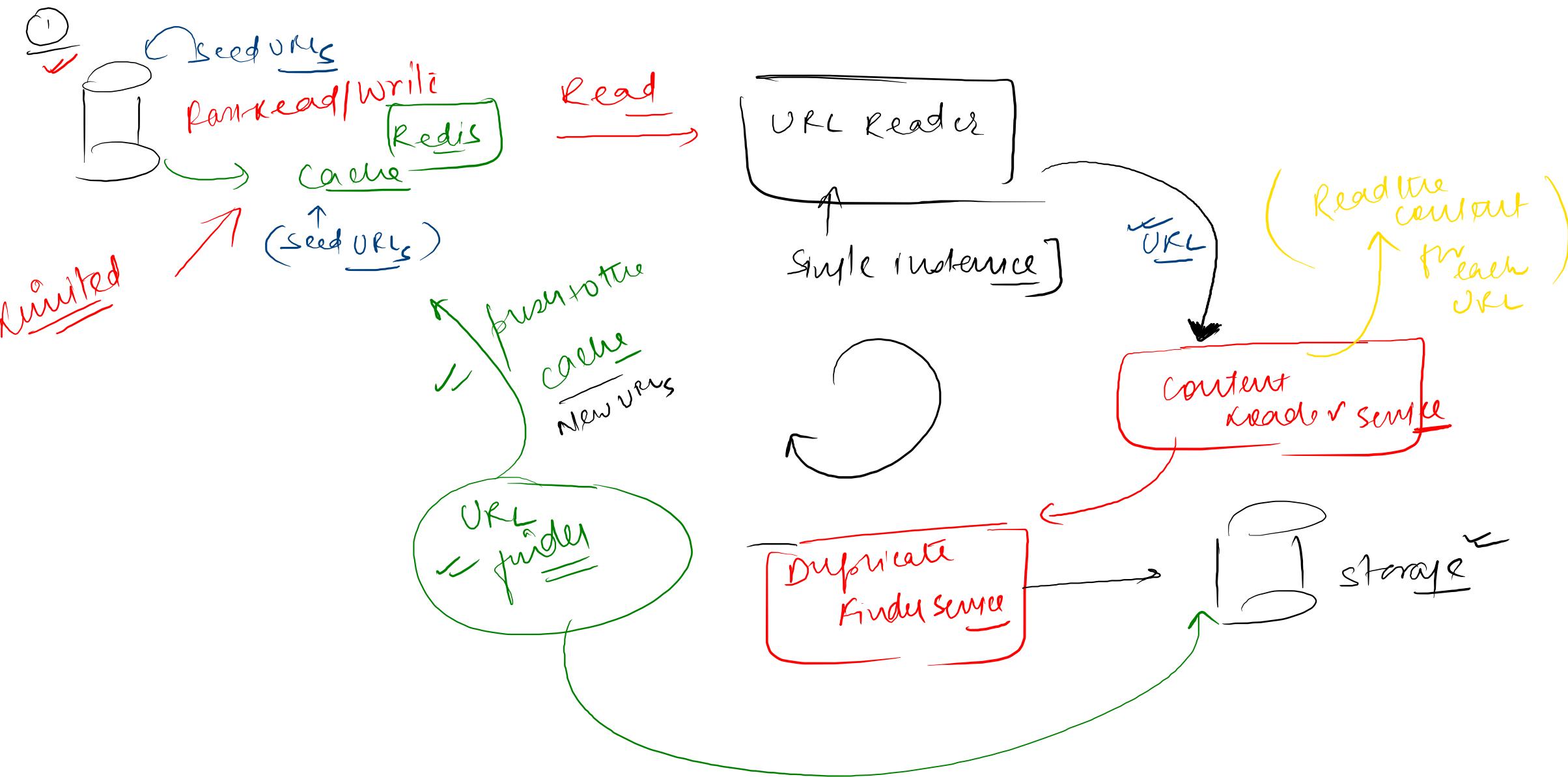
←

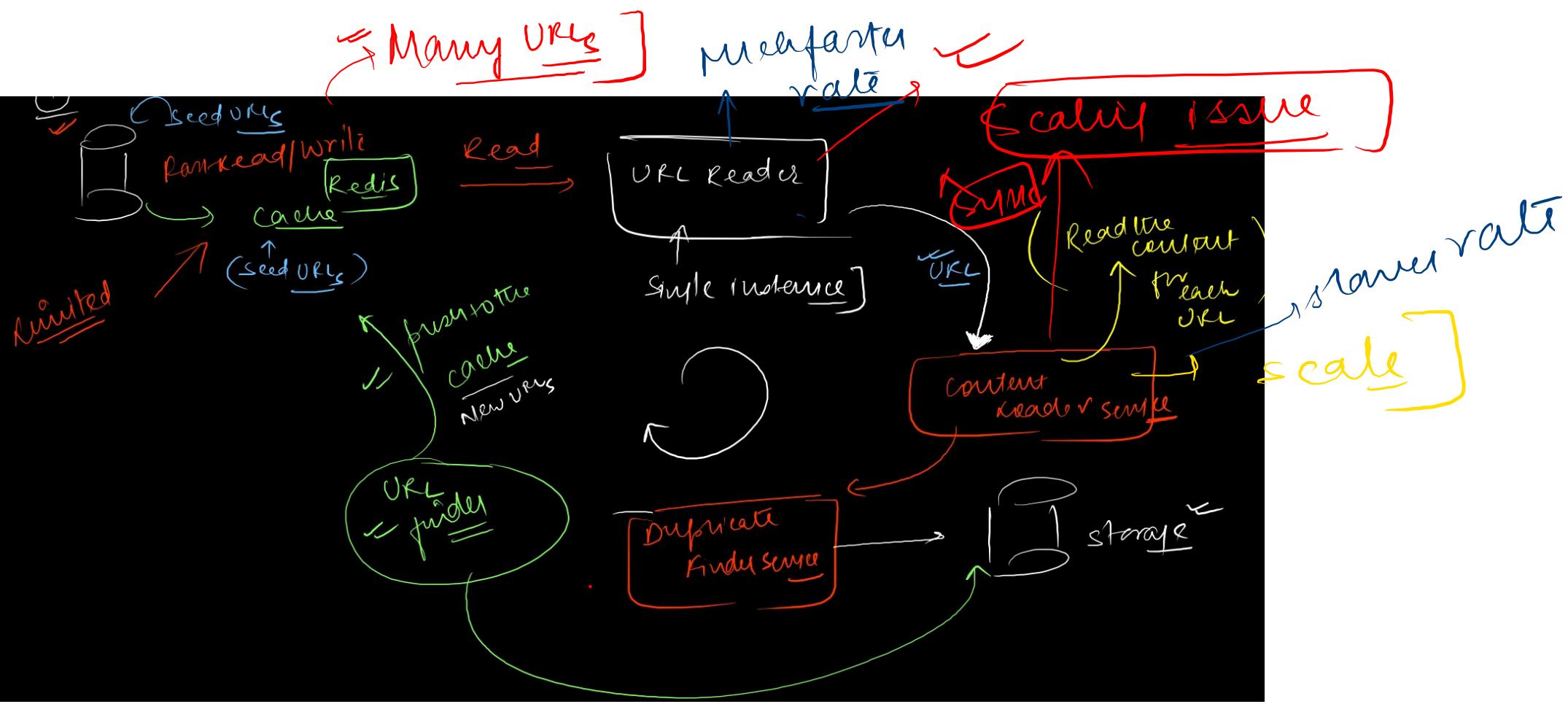
←

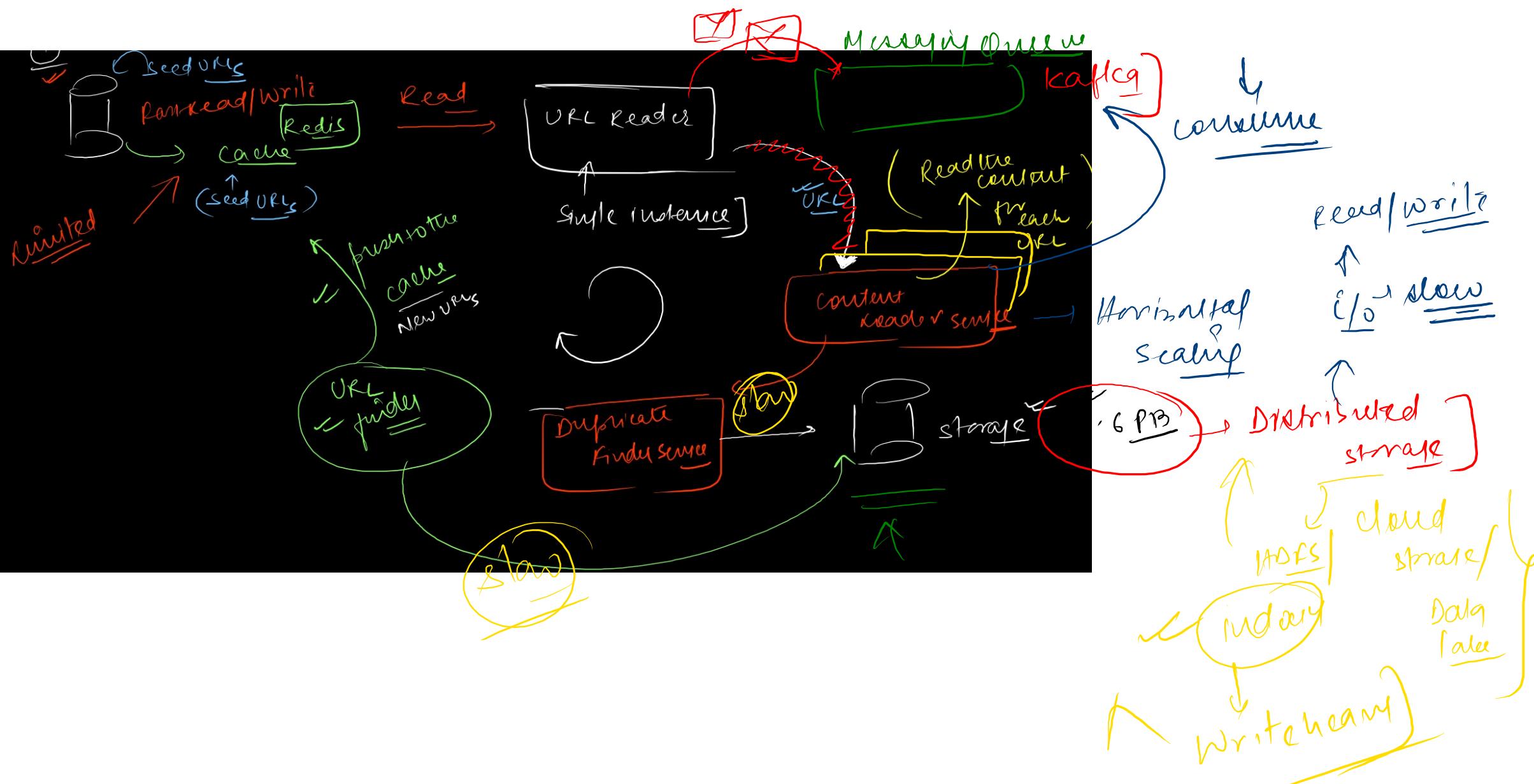
web crawler

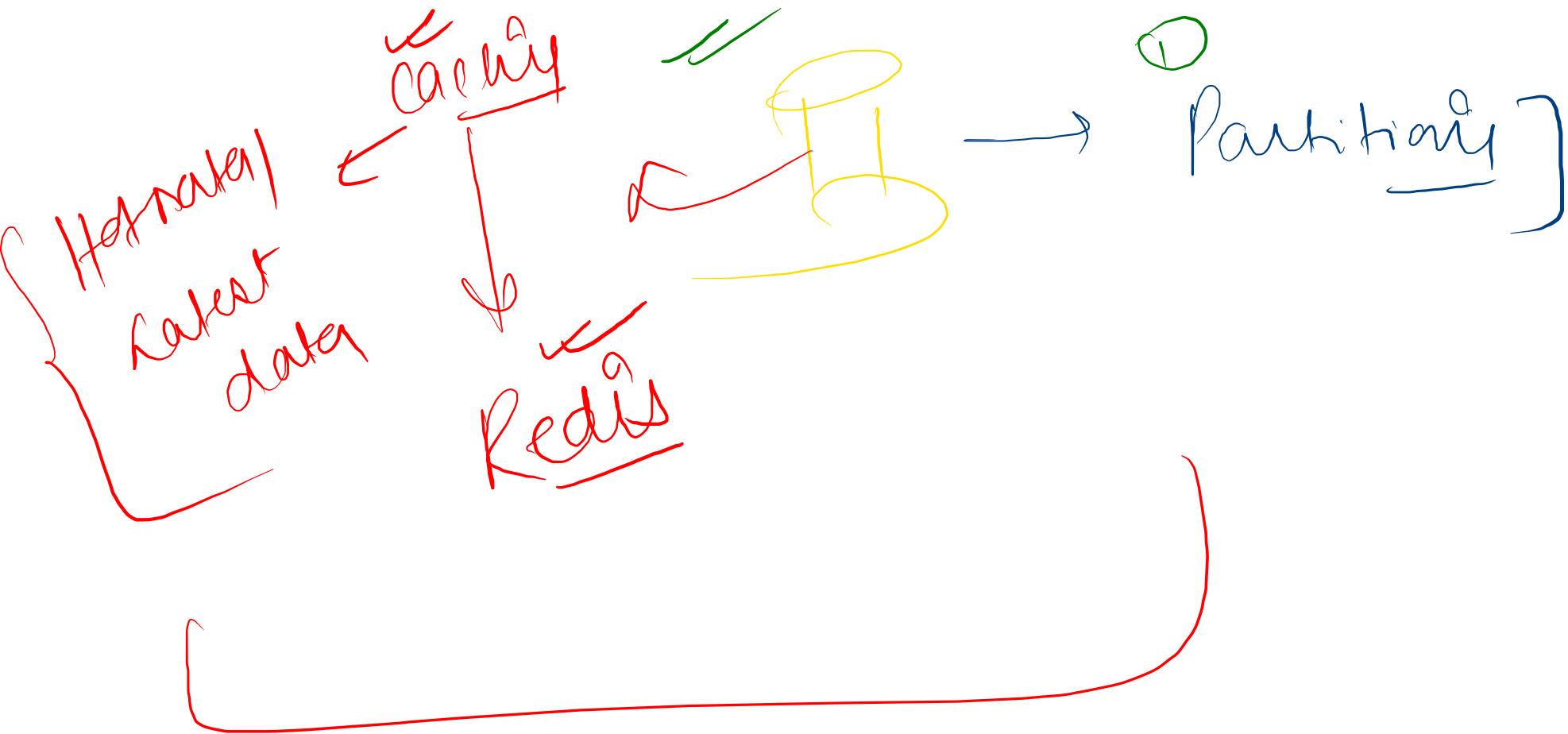
←  
BFS











T4T

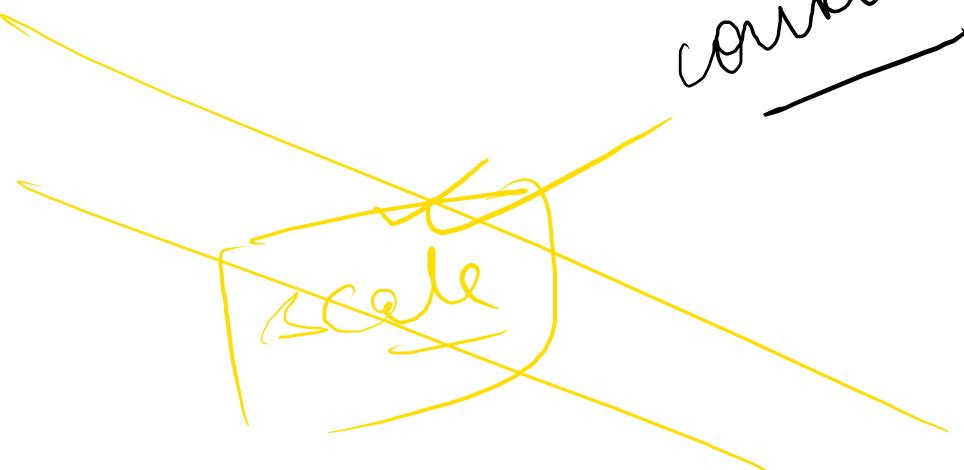
UFL Reader

Kafka

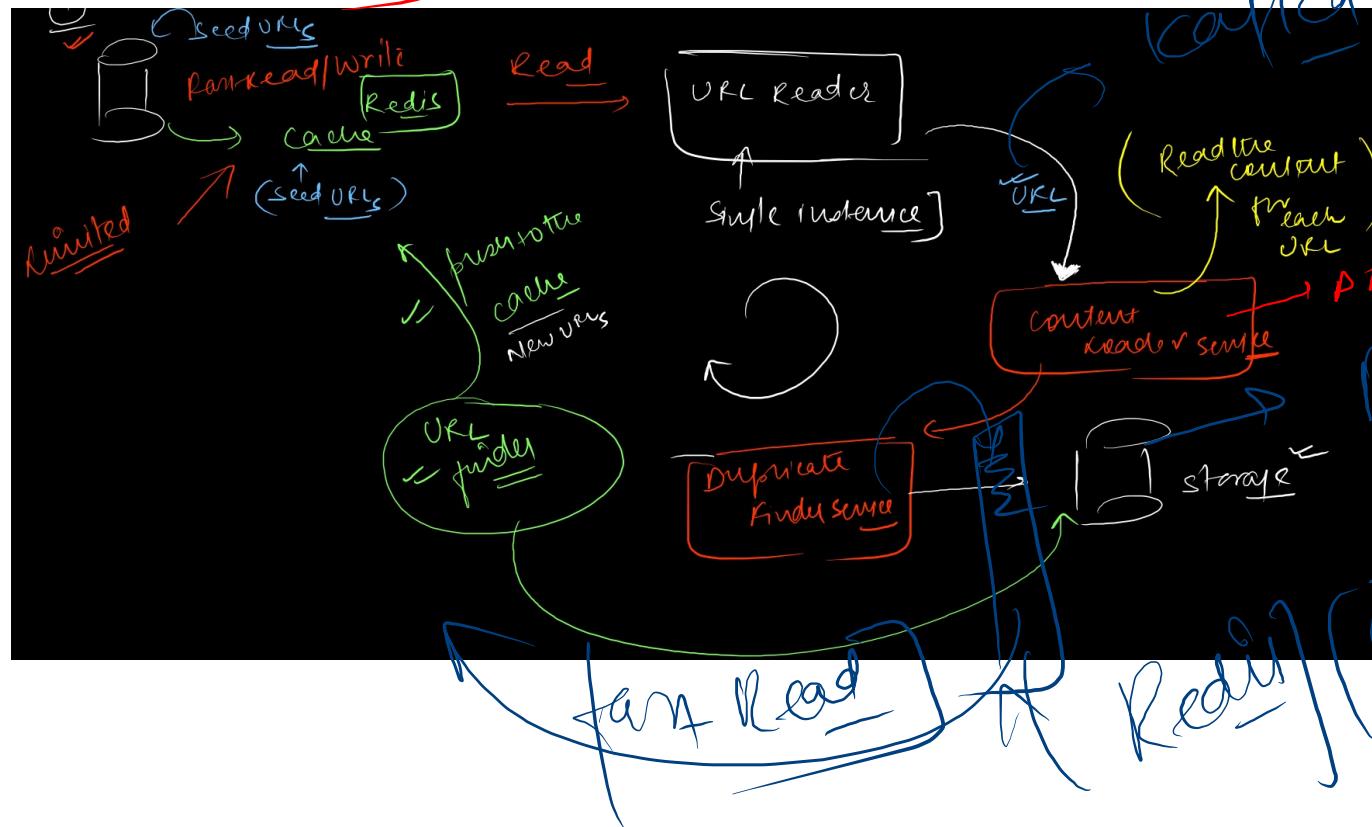
consumed

scale

Content needed  
Service



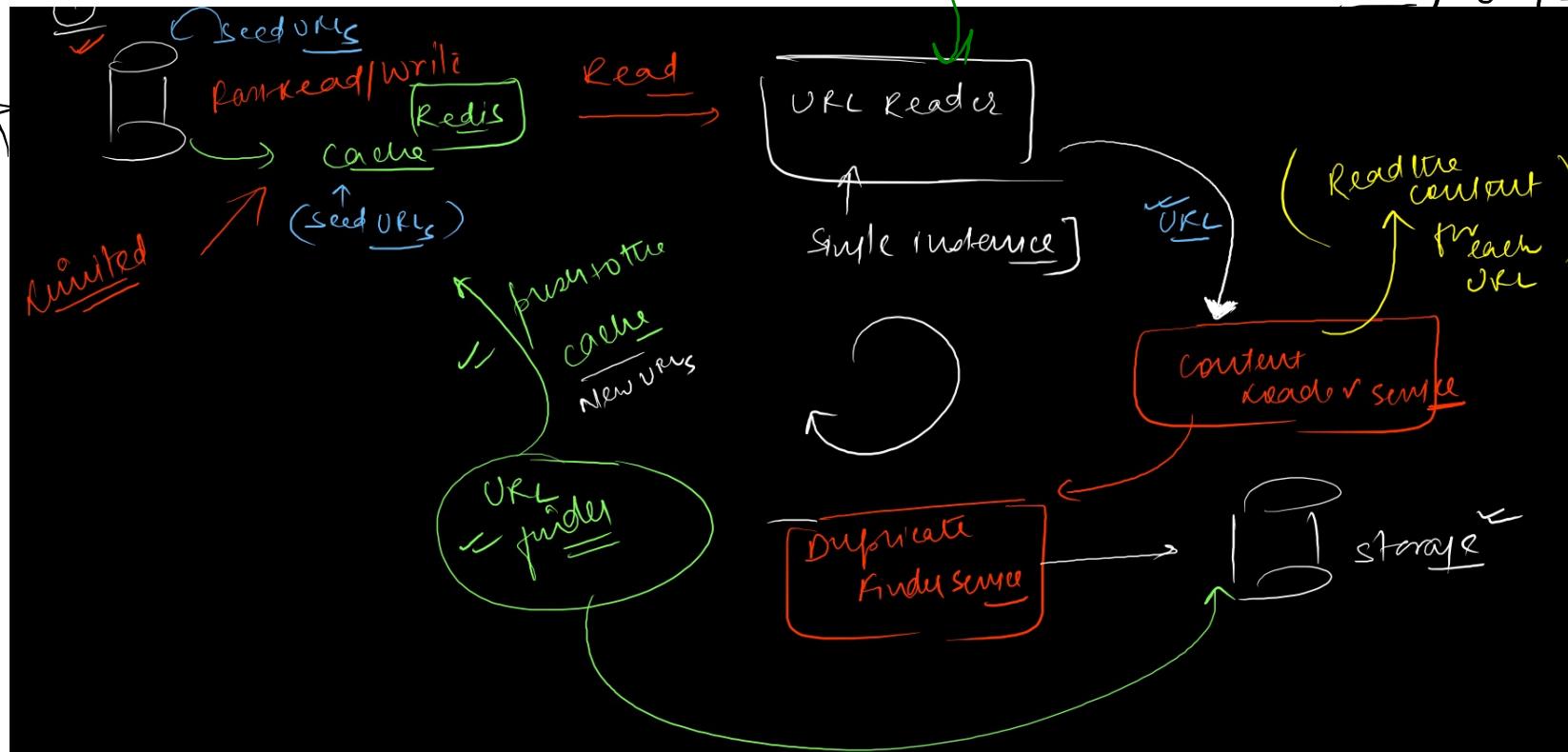
## ✓ [Web Crawler]

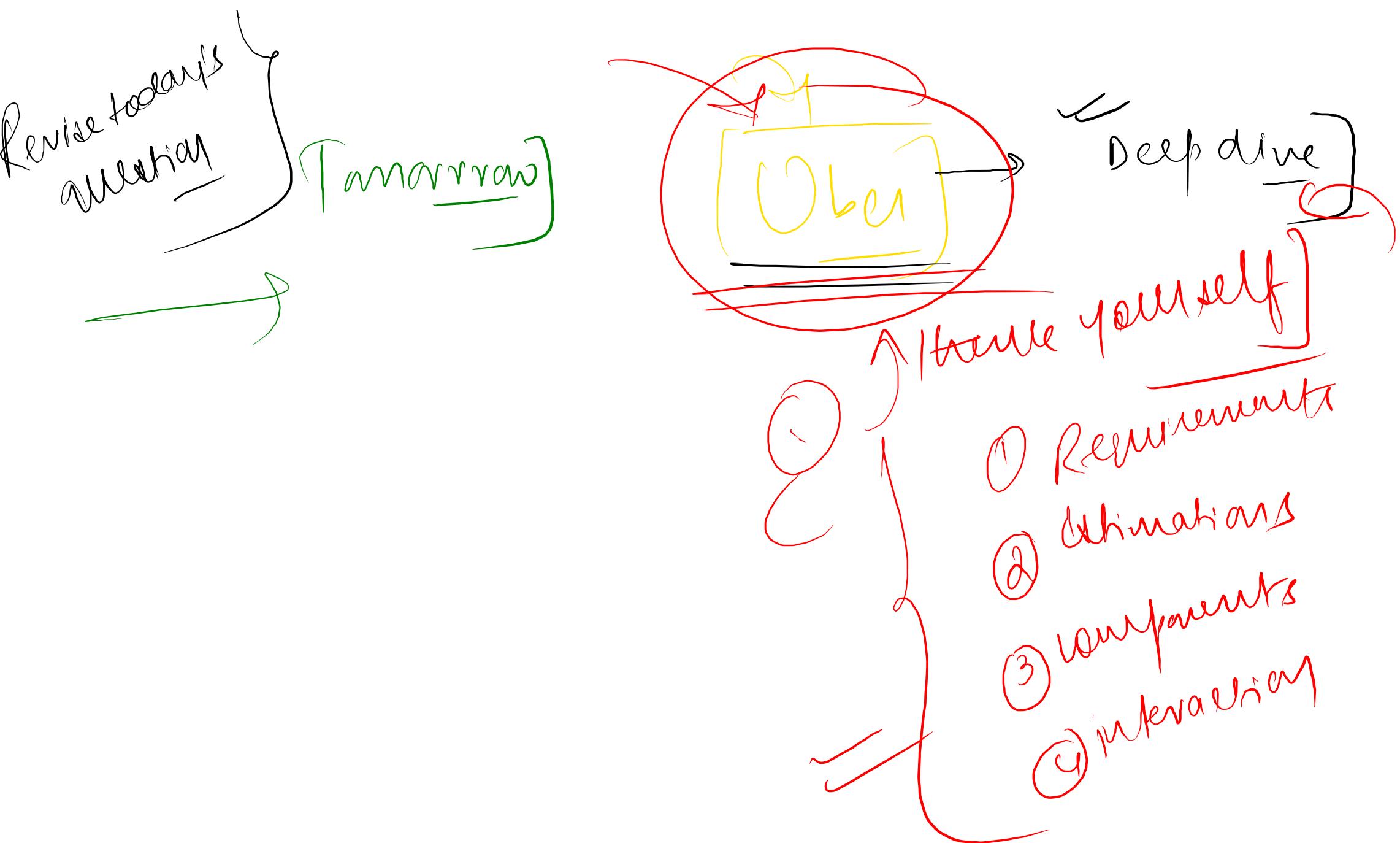


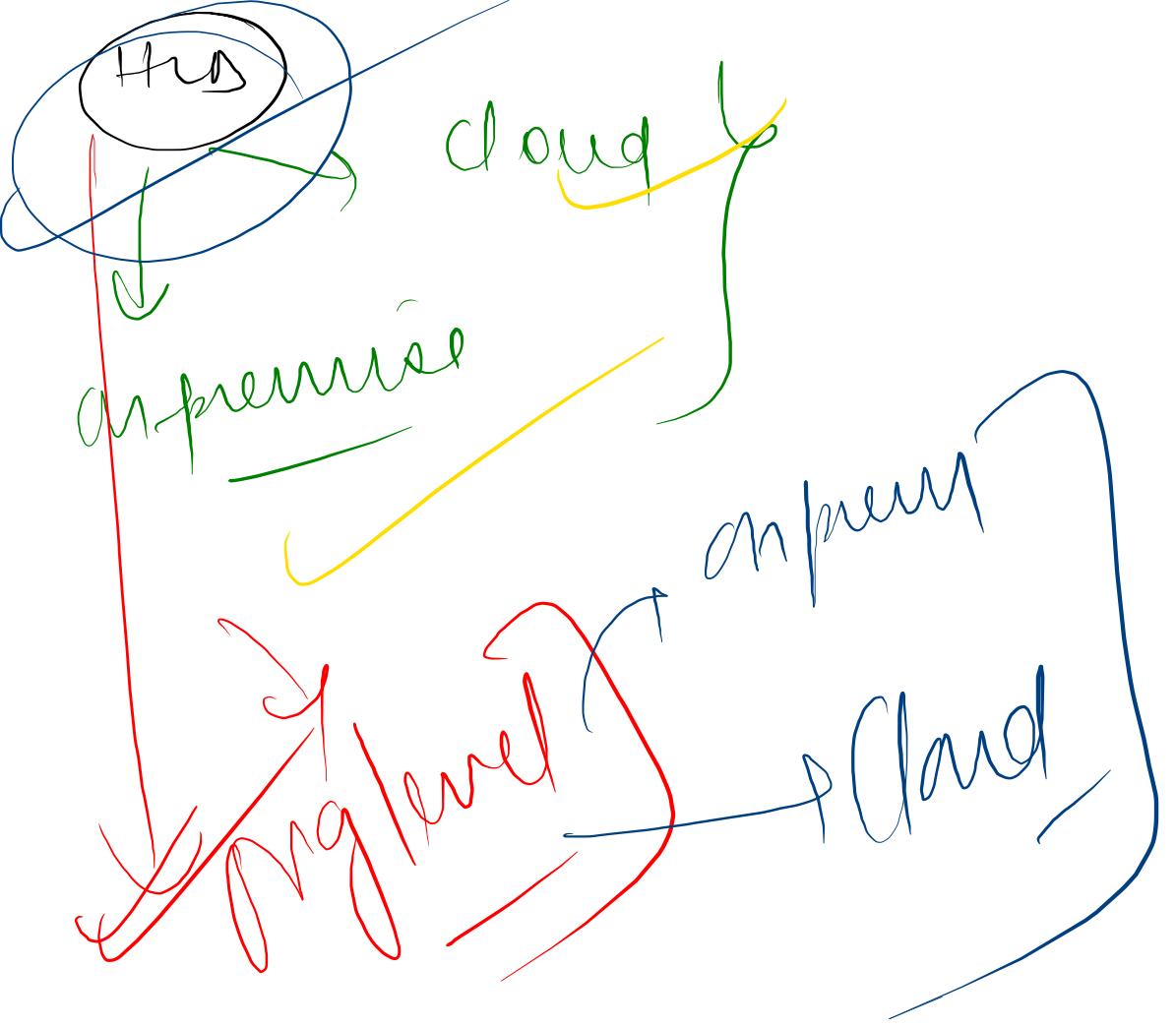
[in b/w]

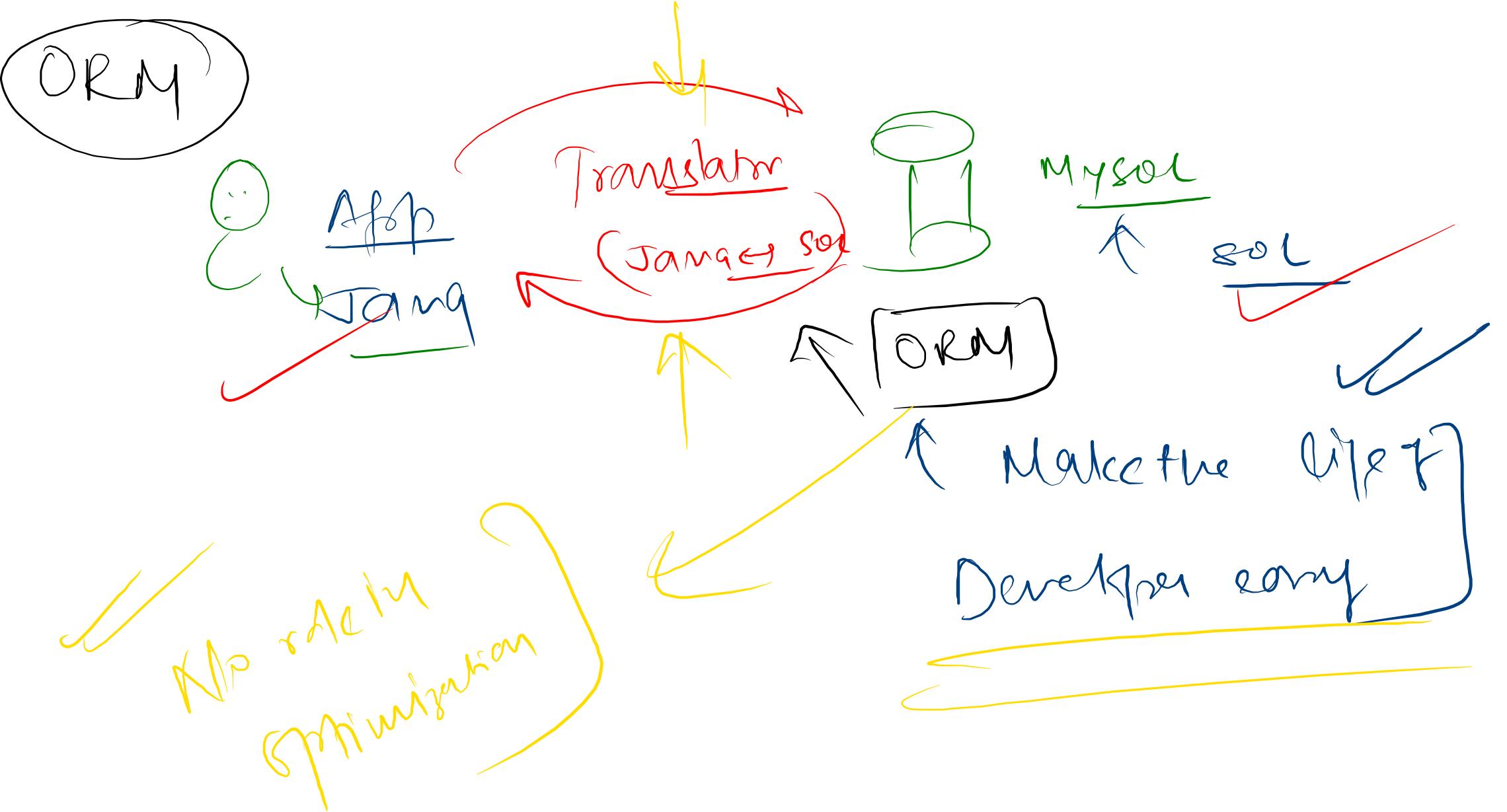
Scope of improvement

- ncache → improve the speed of the system
- Partitioned







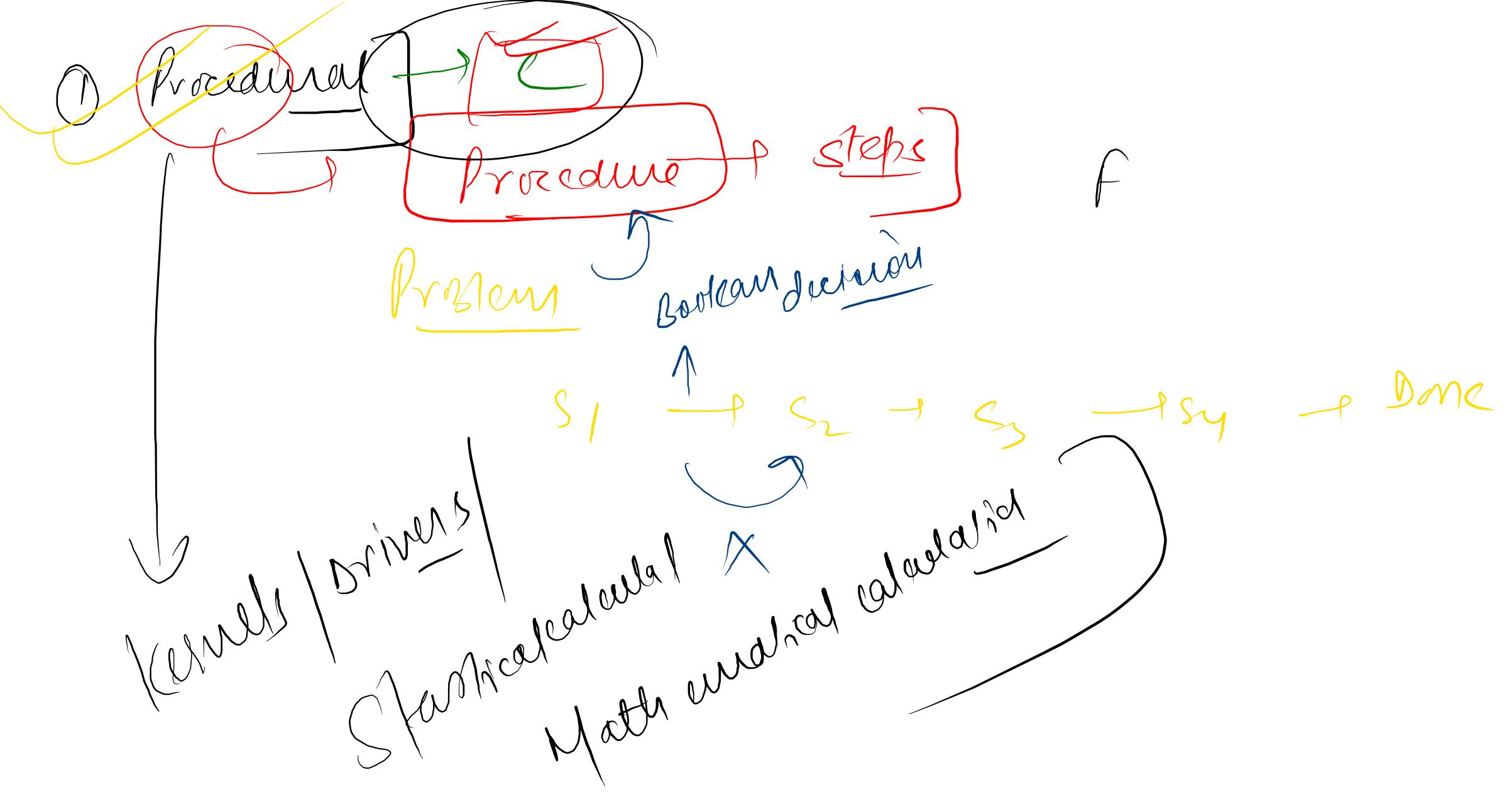


C/C++ / Java / Python / →

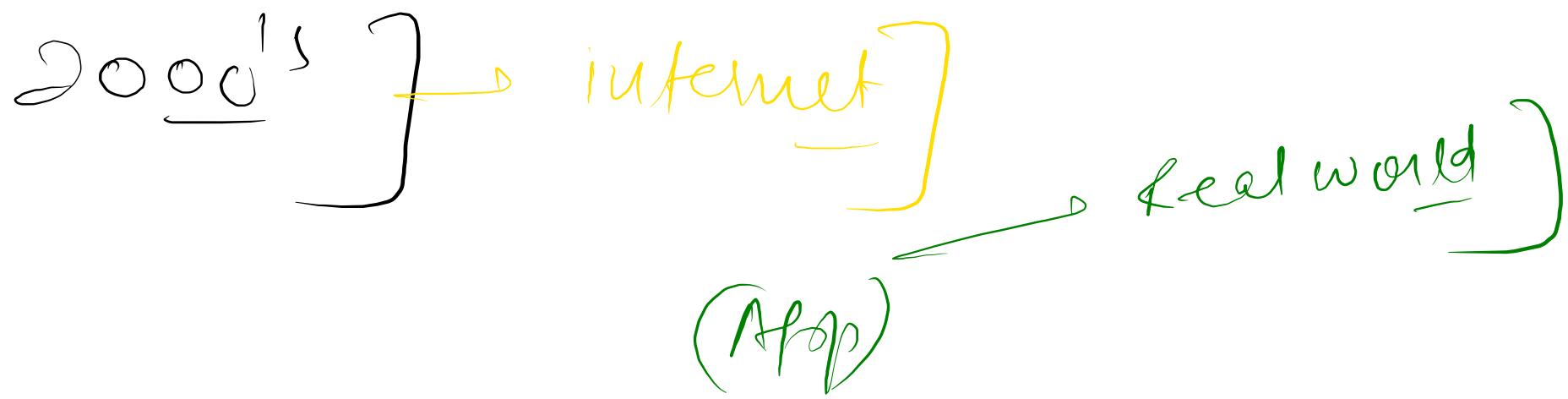
✓ Which to use when

2024

→ Pre-vampy styles / Pandismus ↗

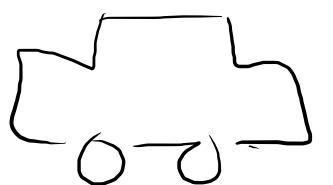


②

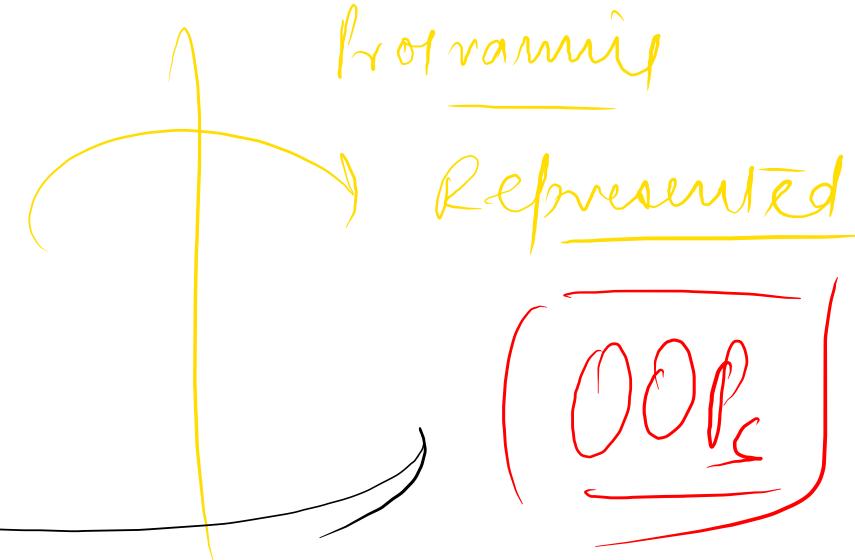


Real world

①



Relationships



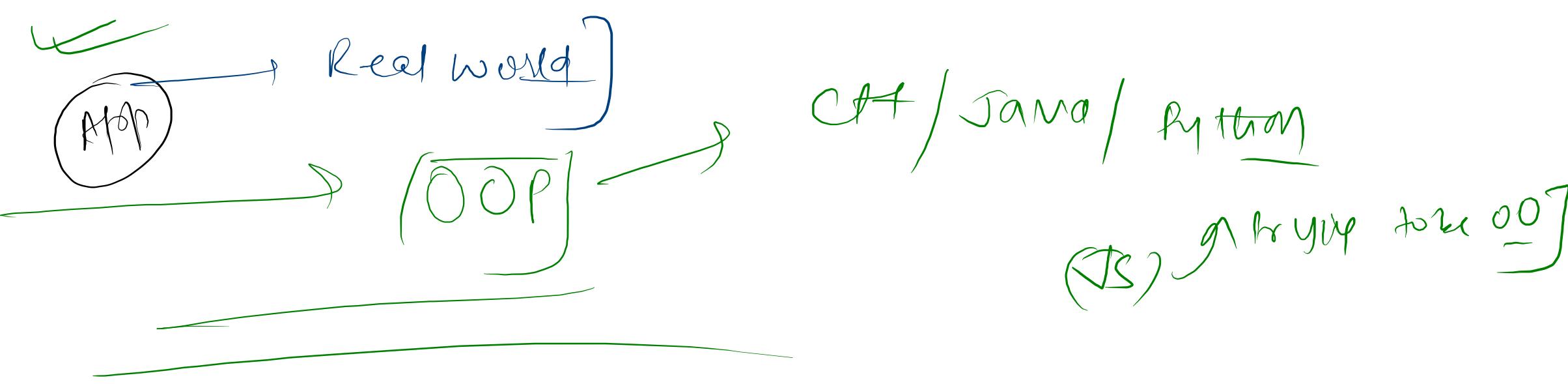
Programming

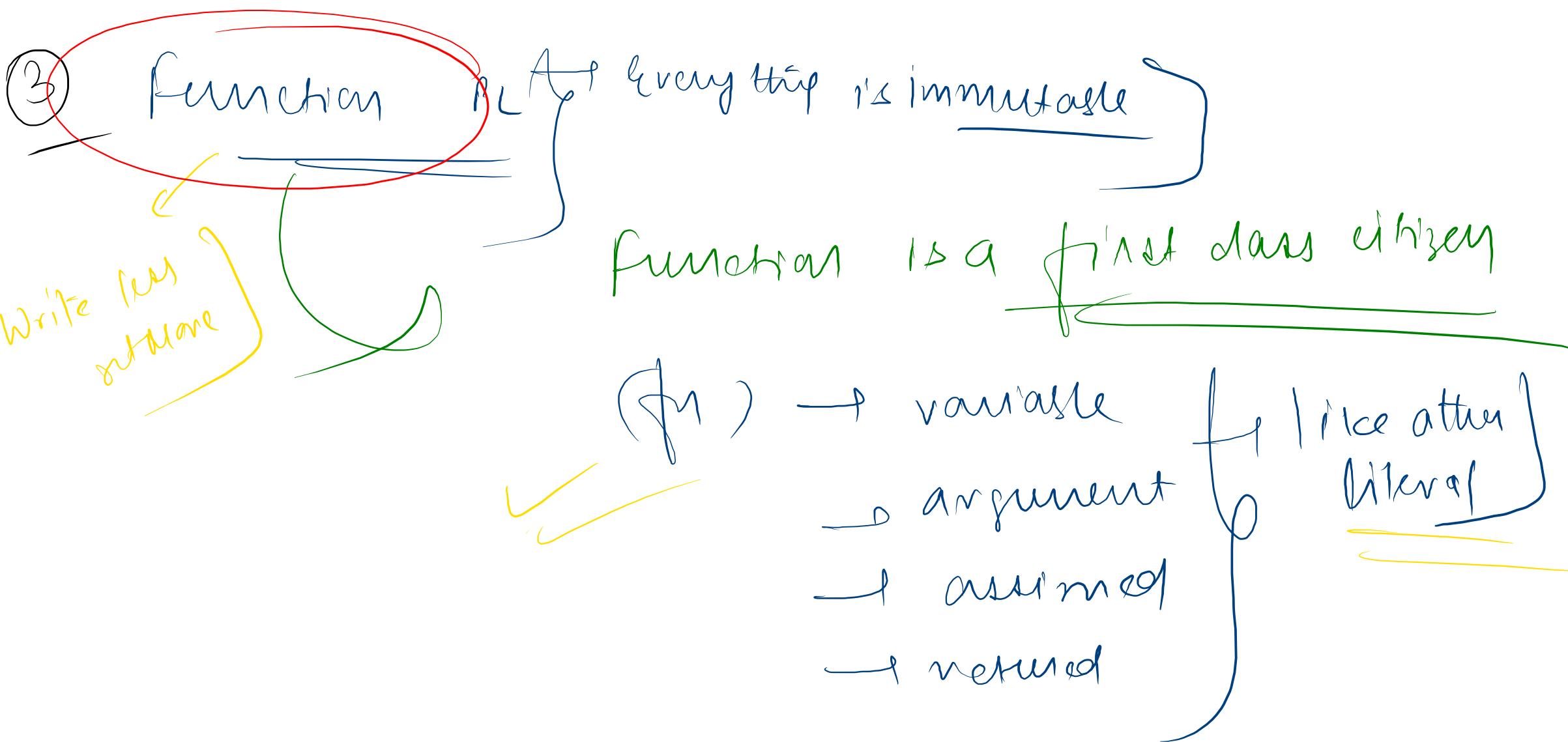
Represented

OOPs

Java

C++





writers → Get Money ] concurrent app

U G M

→ Scala

4

Declarative programming

~~what~~

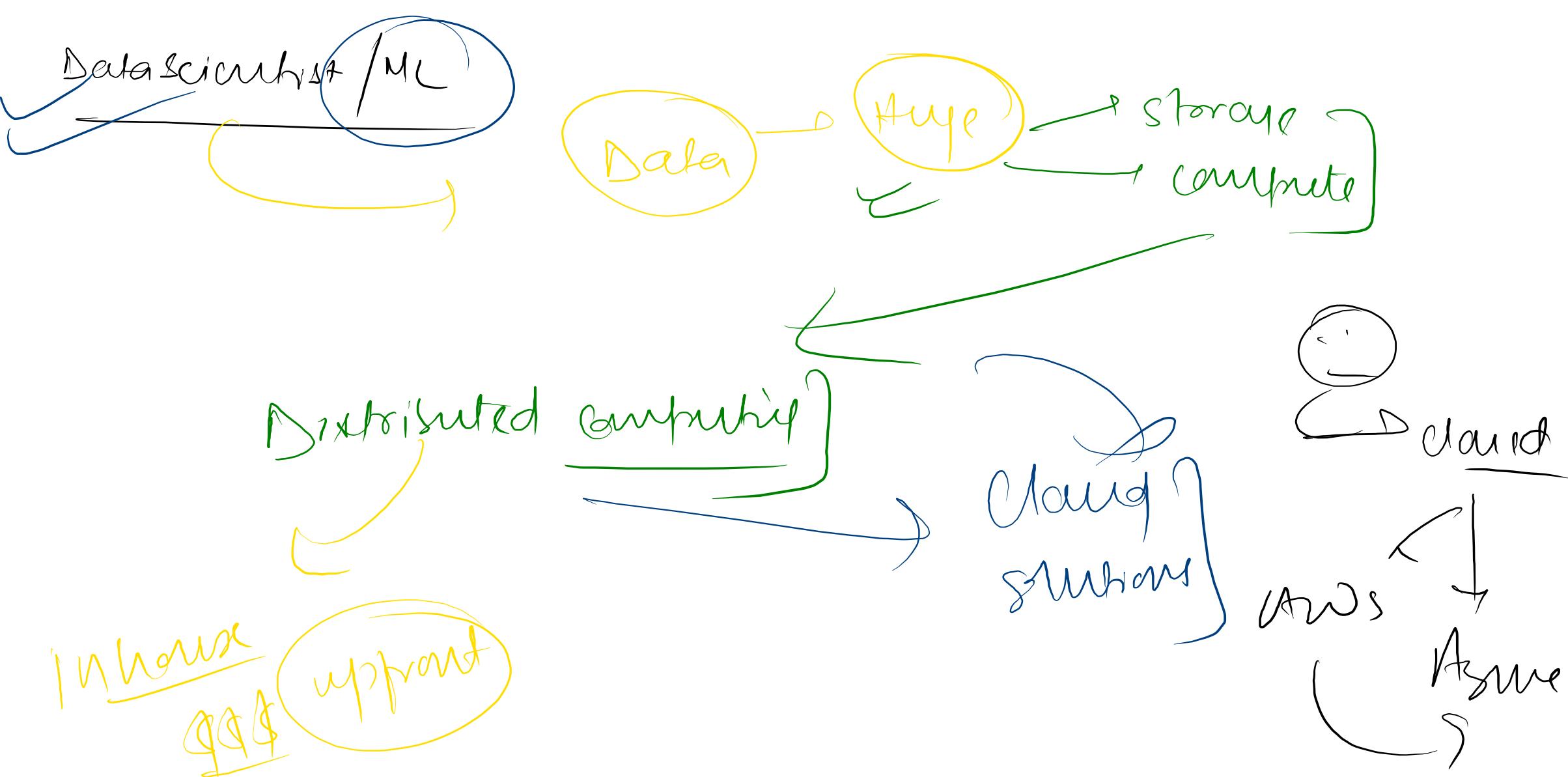
and ~~not how~~

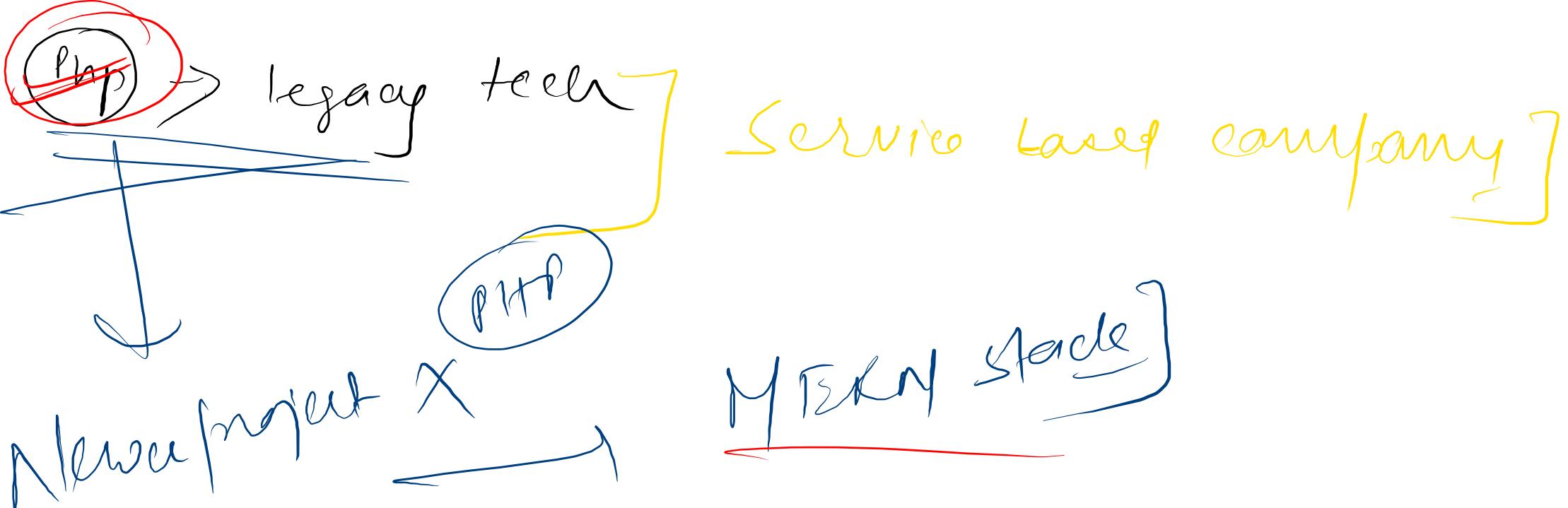
~~what~~

Programs

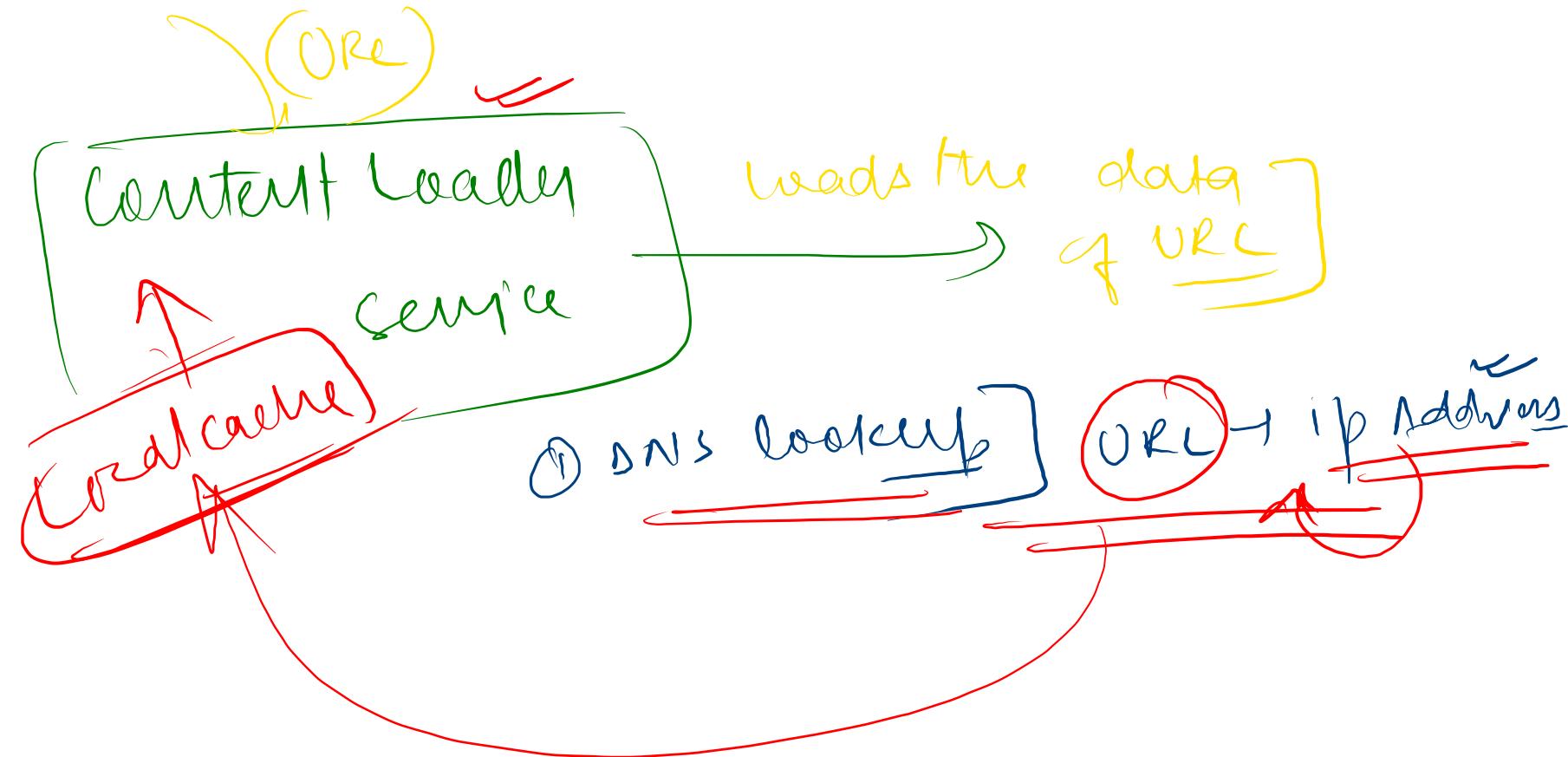
~~what and How (logic)~~

Solved ~~from students;~~





Bigest Player



Caching

speed the read

In memory

✓ cleaned after  
every request

✓ Memcache

Cache

Memcached

