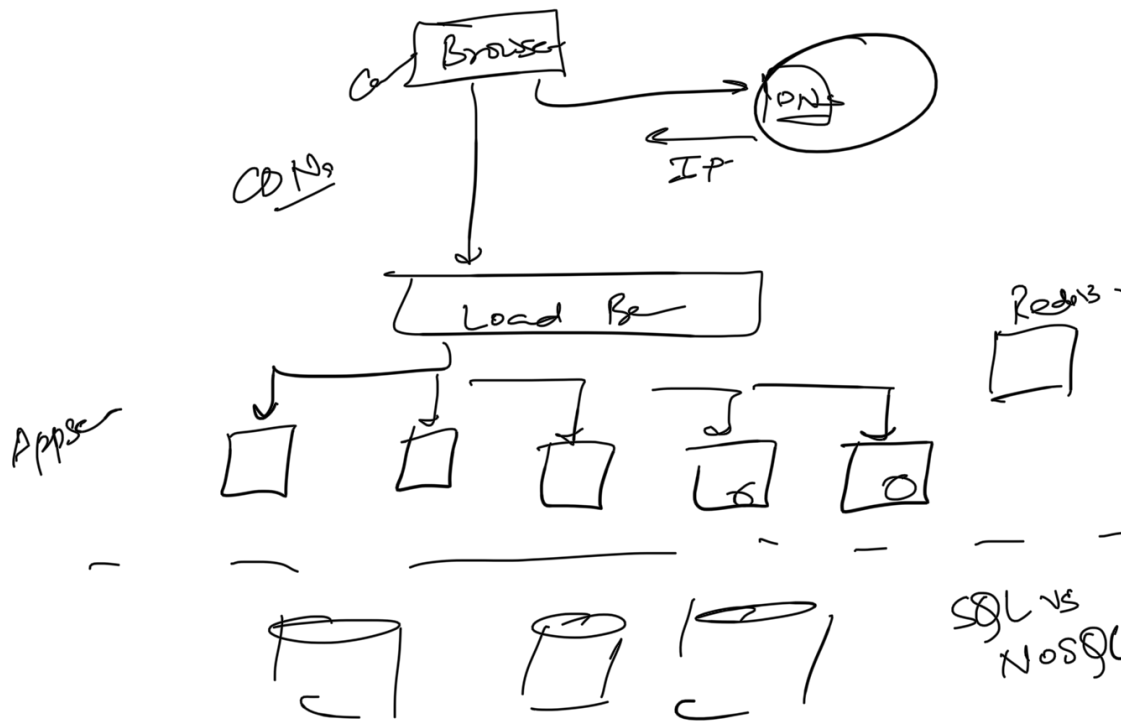


SYSTEM DESIGN

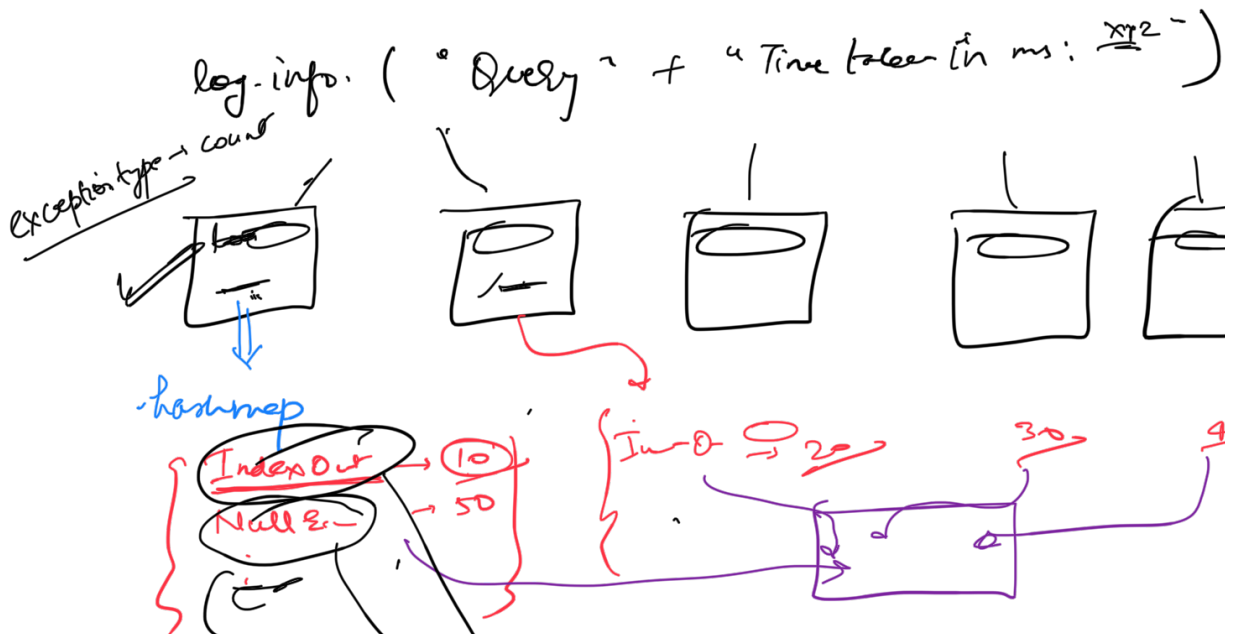


- ① zookeeper , ② ElasticSearch or lucene , ③ Messaging queues
Kafka, Kinesis
- ④ Large file storage
S3, HDFS





log.error("exception.type" + "exception.~~message~~ +
current-ti - ")



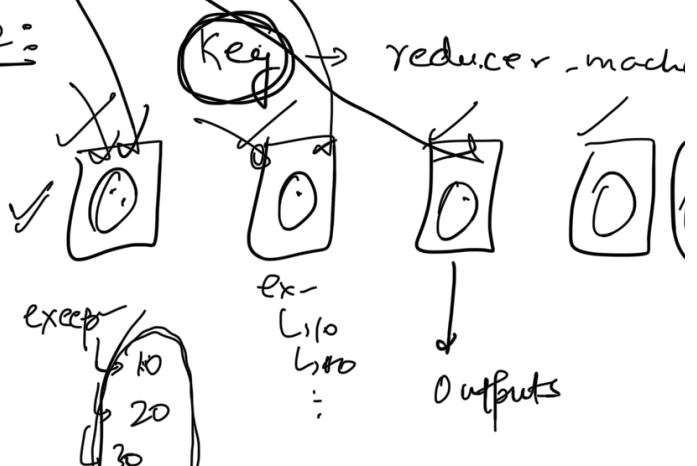
STEP 1: On every machine, go log entry by entry

entry \rightarrow (key, value)

STEP 2:

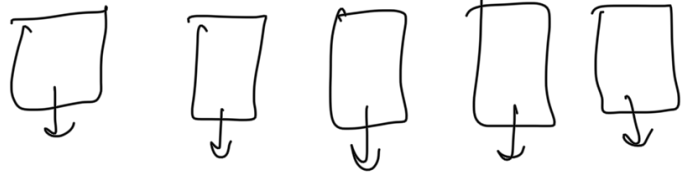
Reducer Machine

key % N
sort the key
key %



STEP 2: Merge the values for common k

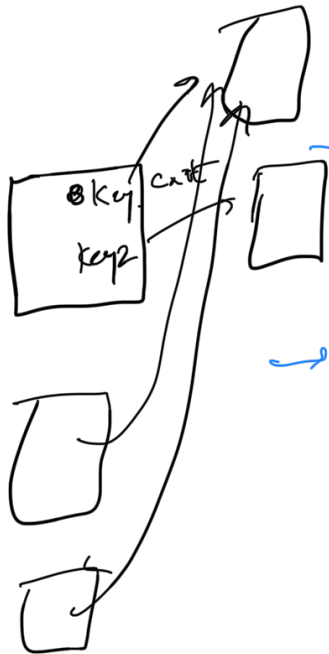
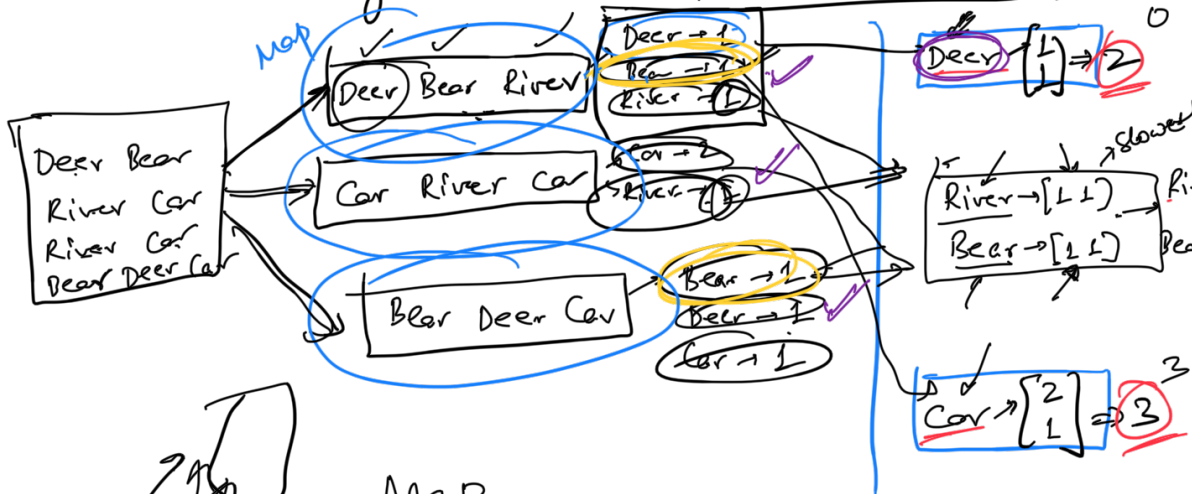
* Exception key \rightarrow total-count



C[0] % 2

Google \rightarrow Map Reduce

4 units
2 units
0



Map

entry \rightarrow key, value, aggregated
Deer \rightarrow Deer (1), sum

Shuffle :

Assign key to reducer node
 \hookrightarrow Independently on map.

Reduce :

list of values

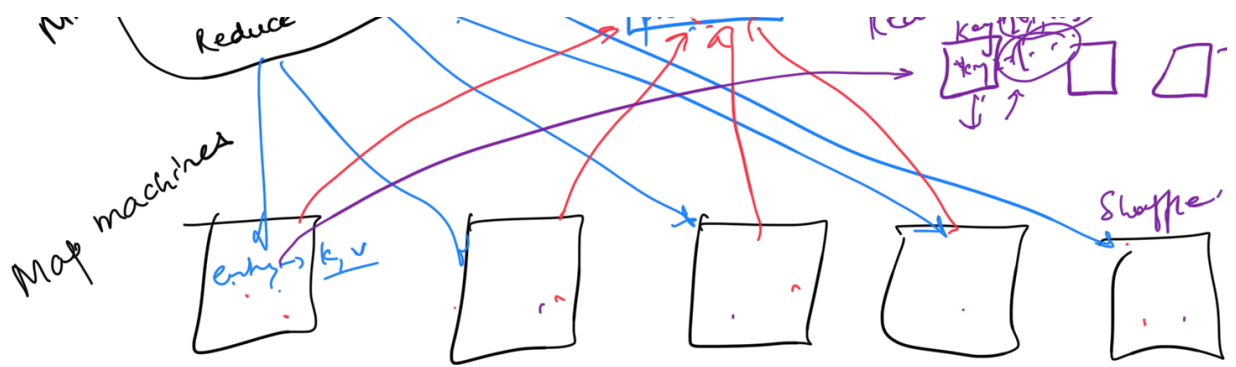
sum, avg, min, max...

ArrayList (values)



key \rightarrow ArrayList (val)

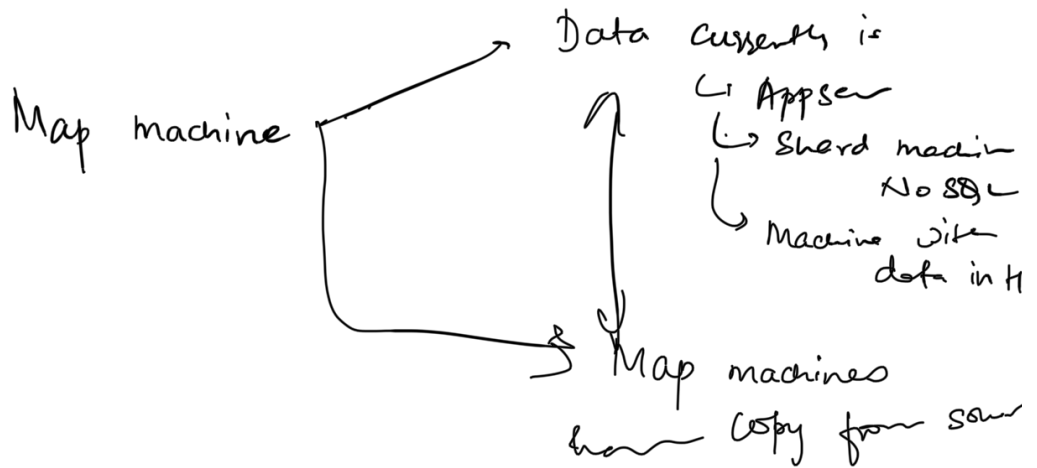
these?



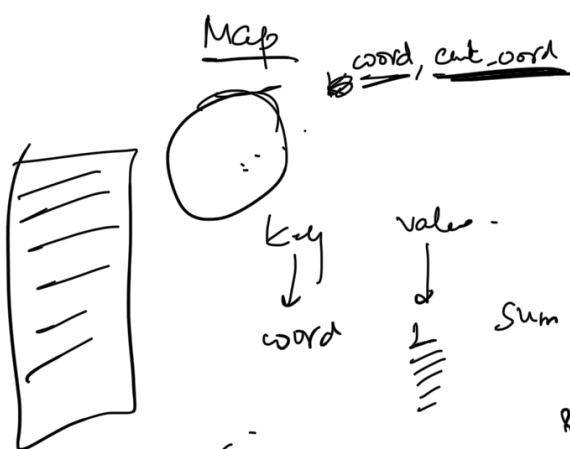
Phase 1: Map completed

Phase 2: Shuffle/sort

Phase 3: Reduce (output is written)



key \rightarrow final-count



Reduce
word \rightarrow [...]
Sum the values

word \rightarrow total-count



A, B, C, D, E

(last, Ar)

$(A, B) \rightarrow A \rightarrow B$

✓ Map

$(X, Y) \rightarrow [$ common for

✓ Reduce

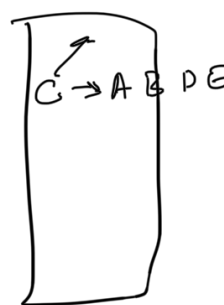
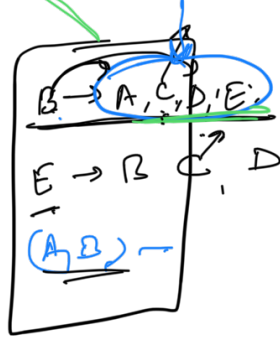
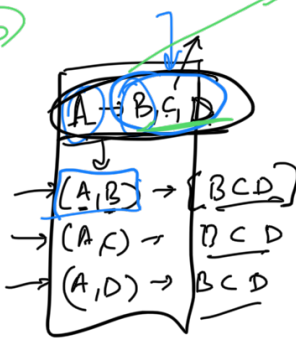
entry

$X \rightarrow [x_1, x_2, x_3, x_4]$

$(\min(x, x_1), \max(x, x_1)) \rightarrow [x_1, x_2, x_3, x_4]$

$(\min(x, x_2), \max(x, x_2)) \rightarrow [x_1, x_2, x_3, x_4]$

(A, B)



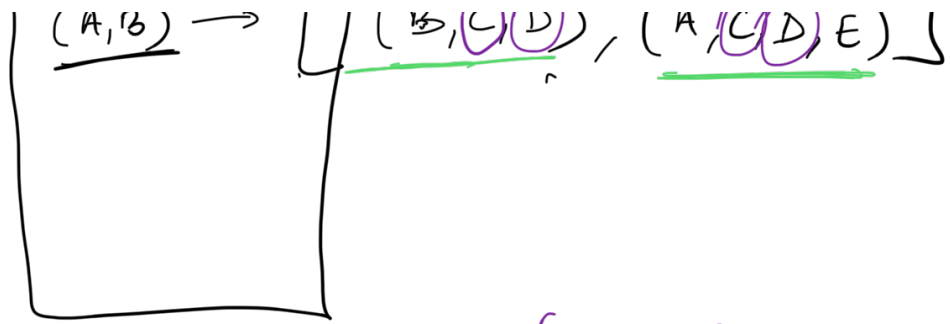
$(A, B) \rightarrow B C D$
 $(A, C) \rightarrow B C D$
 $(A, D) \rightarrow B C D$

$(A, B) \rightarrow A C D E$
 $(B, C) \rightarrow A C D E$
 $(B, D) \rightarrow A C D E$
 $(B, E) \rightarrow A C D E$
 $(C, E) \rightarrow B C D$
 $(D, E) \rightarrow B C D$

$(A, C) \rightarrow A B D E$
 $(B, C) \rightarrow A B D E$
 $(C, D) \rightarrow "$
 $(C, E) \rightarrow "$

$(A, B) \rightarrow$
 $(C, D) \rightarrow$





Map
 $A \rightarrow [B, C, D]$

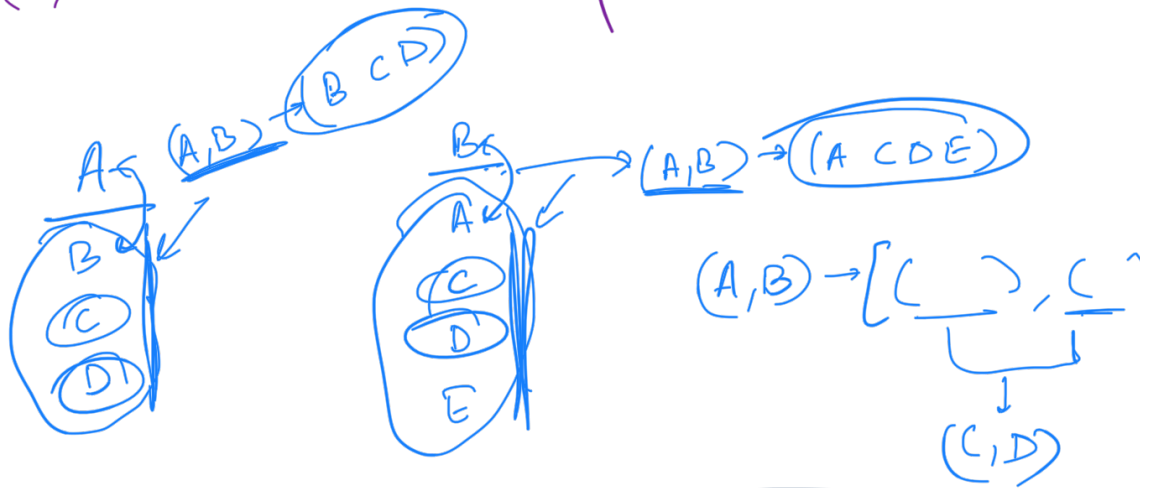
$(A, B) \rightarrow B, C, D$
 $(A, C) \rightarrow D, C, B$

Reduce

$(A, B) \rightarrow [X_1, X_2]$

intersection(X_1

$(A, B) \rightarrow (C, D)$



$(A, B) \rightarrow (C, D)$

Distributed computing \rightarrow Parallel



$\frac{\# \text{ of friend pairs}}{1 \text{ Trillion}} \rightarrow 2 \text{ machines}$

[] []

$\uparrow \quad \uparrow$

1B users
 \downarrow 1000 friends

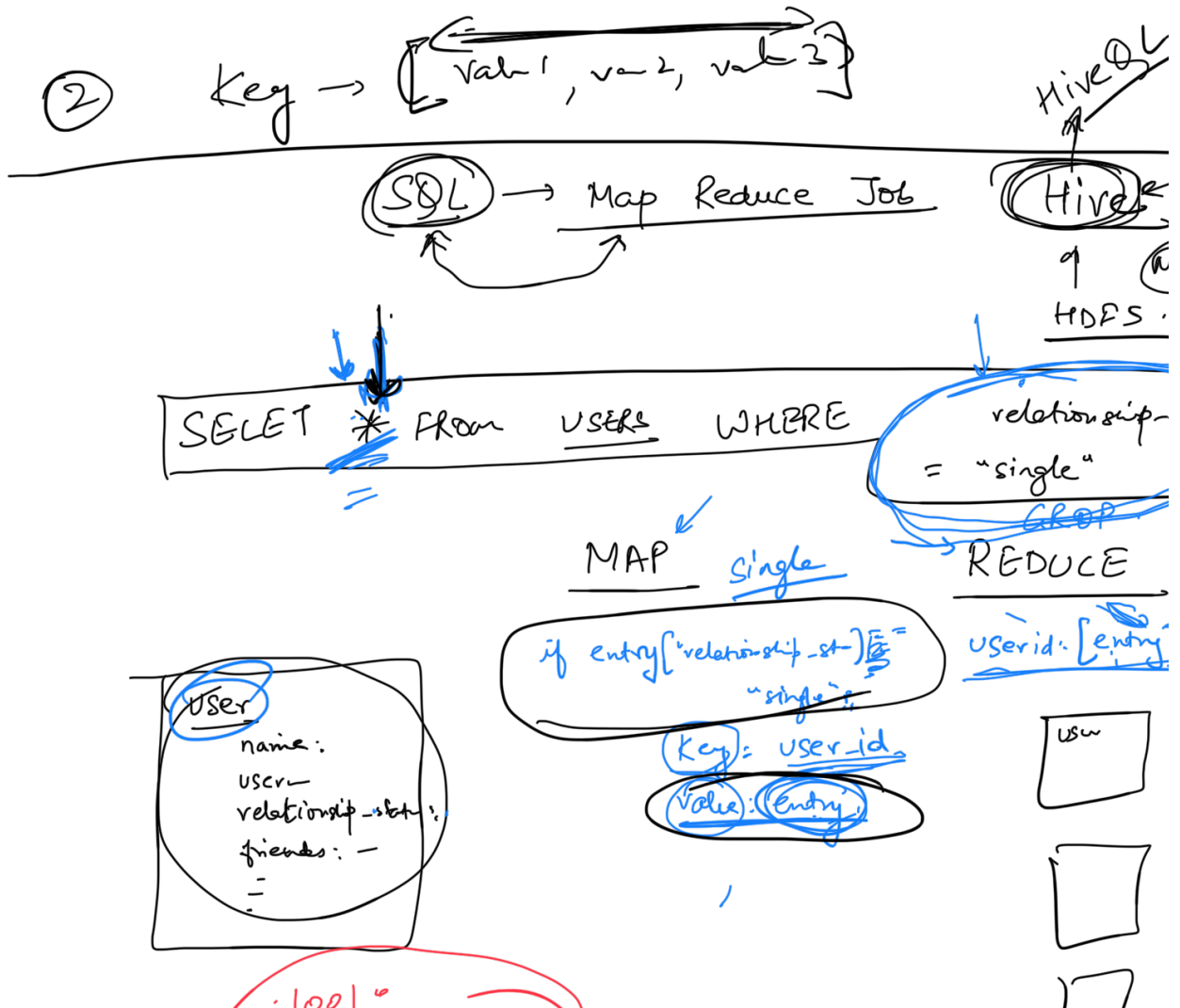
1 Trillion x
 Time takes to
 talk to 2
 me ~~me~~
 Time take to

Map: 1000 * 1000 \leftarrow < 1s.
 similar time \leftarrow < 1s

Reduce: 1 million key \rightarrow $\left[\begin{matrix} L1 \\ \uparrow 1000 \end{matrix} , \begin{matrix} L2 \\ \uparrow 1000 \end{matrix} \right] \leftarrow \begin{matrix} O(1000) \\ \log f \end{matrix}$
 \rightarrow 35 - 10s
10⁹ iterations.

① Entry \rightarrow key, value

② Key \rightarrow [val 1, val 2, val 3]



10-1-0
 name: utkarsh
 username: abcd
 result: single
 for: single

U

user101
 user120
 user121

user1

user1

→ Userid

Userid	site-url
1	1
2	1
3	1

user121

user121

SELECT site-url, count(1) FROM users
 GROUP BY site-url

⇓

MAP
 (user-id, site-url)

(site-url, 1)
 Userid

~~swagat~~ REDUCE
 site-url
 {cnt1, cnt2, cnt3, ...}

Rank	url
1	abc
2	abc
3	sc
4	abcd

agni | jello
 apu | scda
 apu | goya

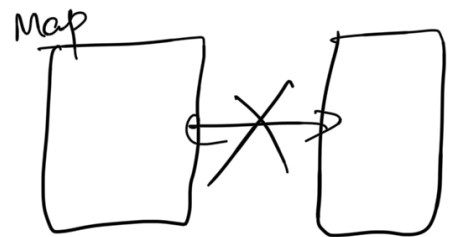
apu

MAP
 (Rank, google.com)
 ↓
 (google.com, 1)

Reduce
 Site → [counts]
 google → [cnt1, cnt2, ...]

Scaler.com

SELECT } MAP
 WHERE }
 GROUP-BY } REDUCE
 HAVING } ↗



She x
 s y
 S-k 2
 S= 2,
 !

100GB → 20GB
 ↓
 1000GB 200G

Bangalore

