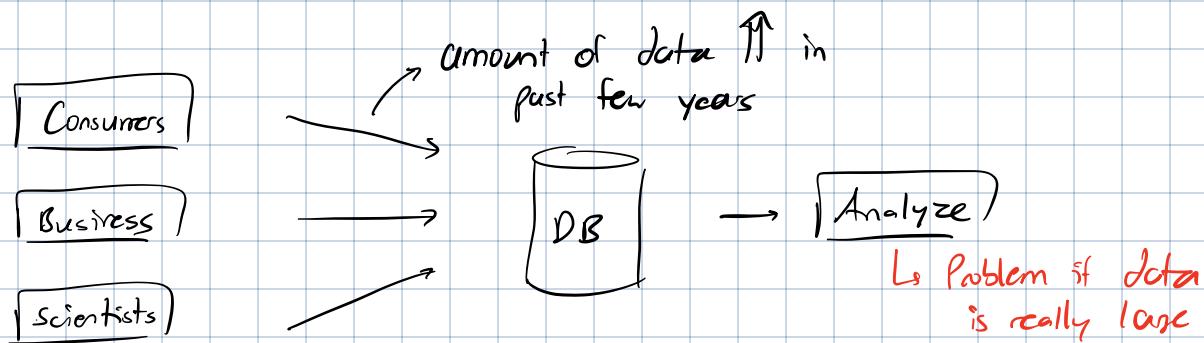


INTRODUCTION

Motivation



Analysis at scale leads to interesting sights

- ↳ Ex:// Target predicting pregnancy, Walmart used big data to design stores
black hole pictures

How to analyze lots of data at scale?

Vertical scaling

- Really expensive HW can do work
- Cap on how much scale to achieve

Horizontal scaling

- Get more computers
- Cheap computers work together
 - ↳ Parallelization
 - ↳ Fault tolerance

focus of course

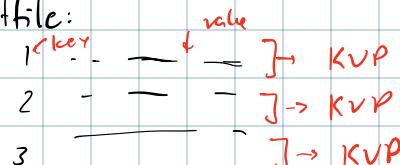
We will use distributed computing frameworks as an abstraction

MAP REDUCE

The Algorithm

Revolves around key-value pairs

- ↳ Ex:// In a textfile:



Functions

Map:

$$(k_1, v_1) \rightarrow \text{list} [(k_2, v_2)]$$

not the same!

Reduce

$$(k_2, \text{list}[v_2]) \rightarrow \text{list} [(k_3, v_3)]$$

all values associated w/ k_2

Basically, extract keys & values from some KVP

Takes all values associated w/ one key \Rightarrow some operation

Ex://

Input:

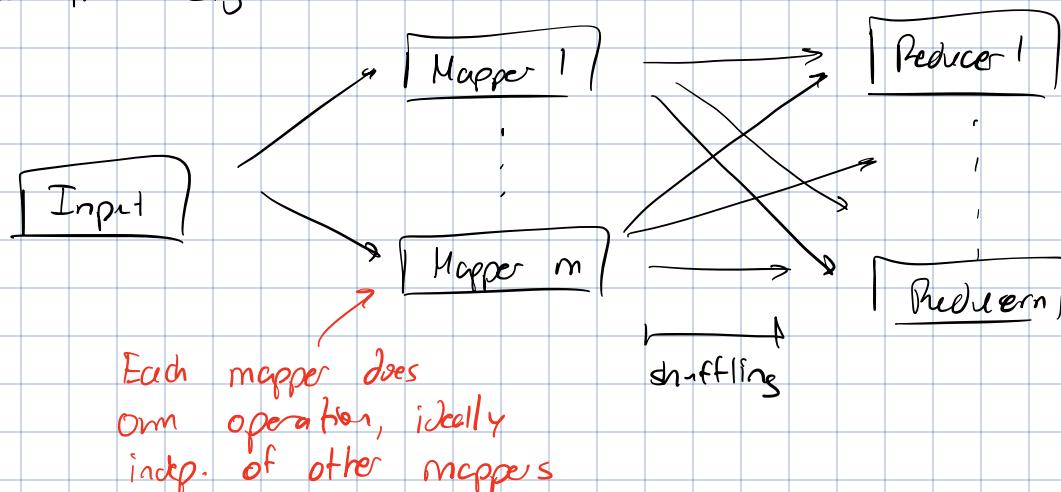
	Output.
0: "abcdaaa"	$\{ \begin{matrix} a: 4 \\ b: 1 \\ c: 1 \\ d: 1 \end{matrix} \}$
1: "bcdadac"	$\{ \begin{matrix} a: 1 \\ b: 1 \\ c: 2 \\ d: 2 \end{matrix} \}$
2: "dabdb"	

Ex:// If aggregation:

$$(a, [4, 1]) \xrightarrow{\text{reduce}} (a, 5)$$
$$(b, [1, 1]) \xrightarrow{\text{reduce}} (b, 2)$$

,

What is the design:



Shuffling: mappers partition particular keys s.t. all keys go to same reducer

↳ Partition: $(k_i, v_i), n \in \mathbb{N} \rightarrow [0, n]$

Assigns k_i to some n ($n = \# \text{ of reducers}$)

Apache Hadoop

Coordinating computing nodes to run MapReduce is hard \rightarrow Hadoop

What does it do:

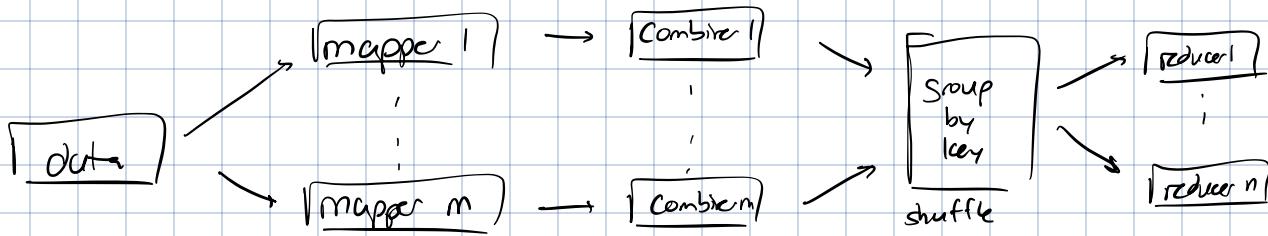
- ① Assign workers to map & reduce tasks
- ② Divides data b/w mappers
- ③ Groups intermediate data \rightarrow sends to right reducers
- ④ Handles errors

Slowest Hadoop op: sending data b/w mappers & reducers

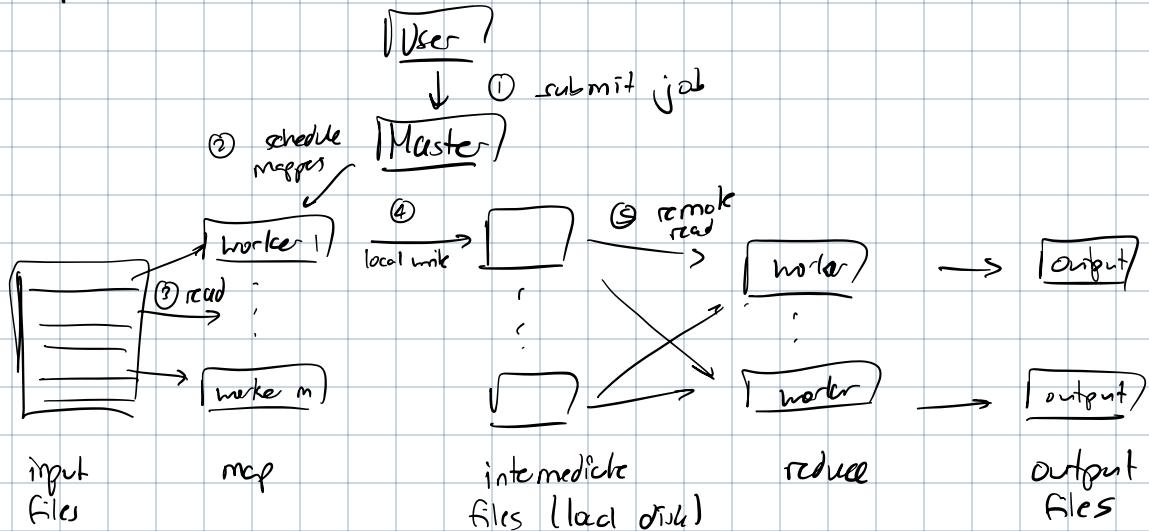
Combiner: $(k_2, \text{List}[v_2]) \rightarrow \text{List}[(k_2, v_2)]$

↳ Take all values associated w/ k_2 & combine. Ideally, $\text{len}(\text{output}) = 1$

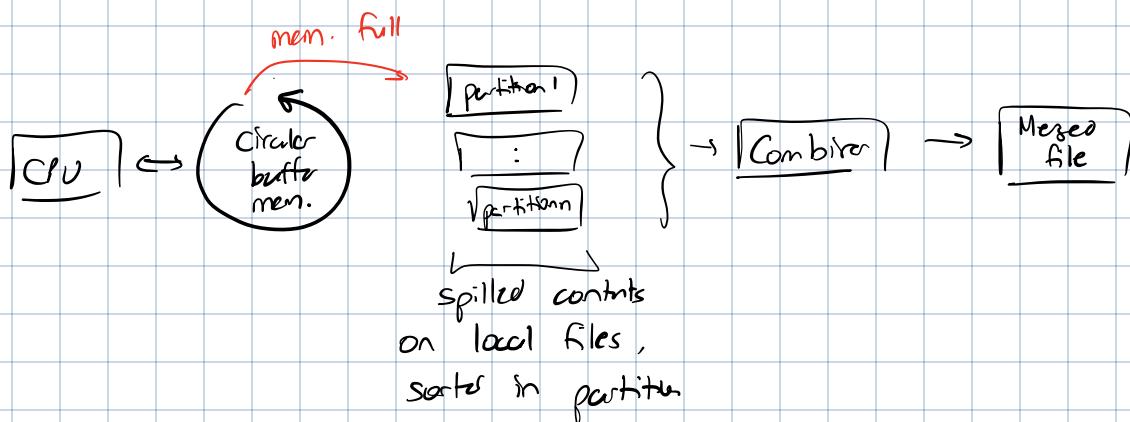
↳ May / may not function like reducer. Doesn't have to run & decided by Hadoop



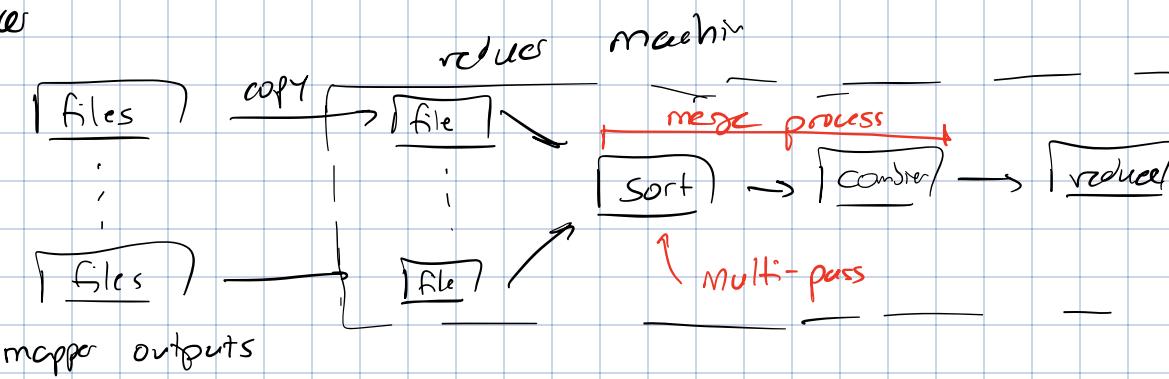
Physically:



Mapper:

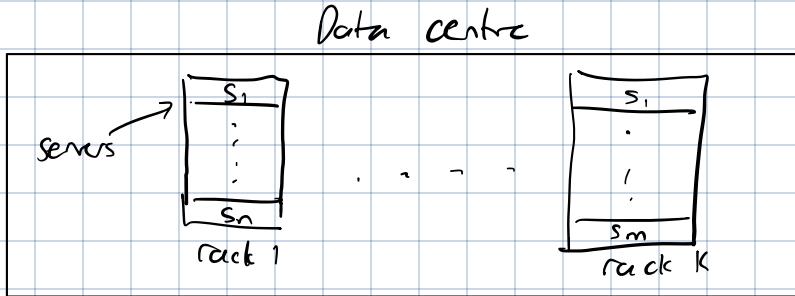


Reducers



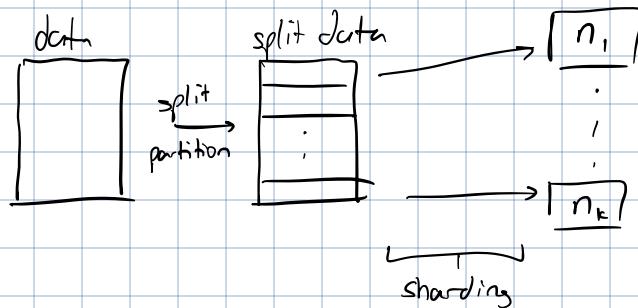
Distributed Storage

Physically:



If we need to perform data ops, locality matters \rightarrow latency \downarrow , bandwidth \uparrow

To store a large amount of data across servers:



- To ensure fault tolerance, store each data partition on multiple nodes (replication)

Hadoop Distributed File System (HDFS)

Can handle 10K nodes, 100M files, 10 PB of data \Rightarrow very large distr. file sys.

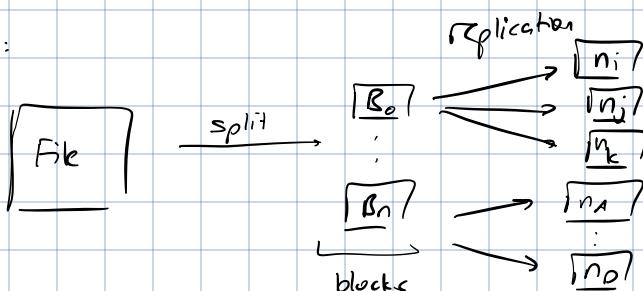
Assumptions:

- ① Commodity hardware: file replication is baked in
- ② Batch processing

Basic principles:

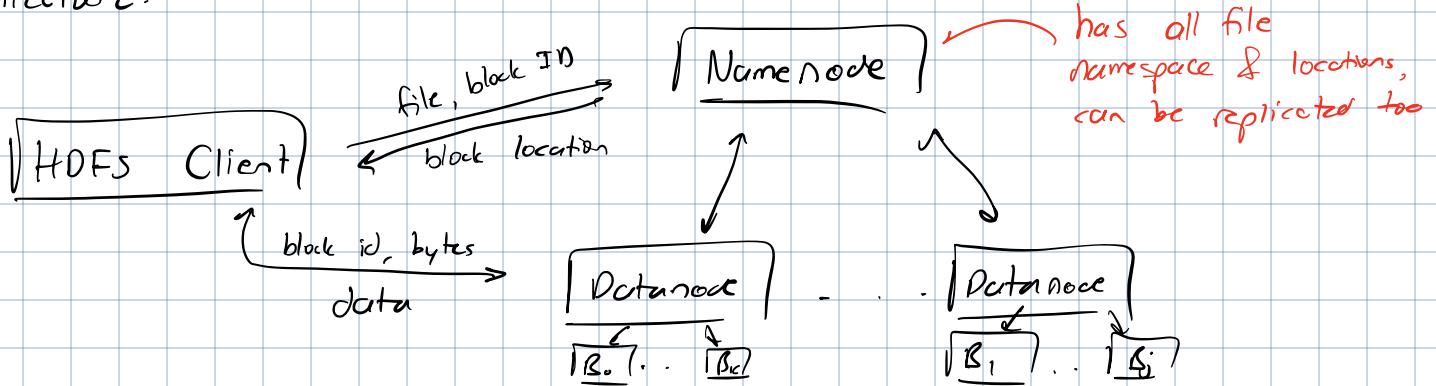
- ① Data consistency: enforces write-once-read-many data model \rightarrow client can only append to files

- ② Blocking:



③ Intelligent client: client knows block location & can access it

Architecture:



Roles of namenode:

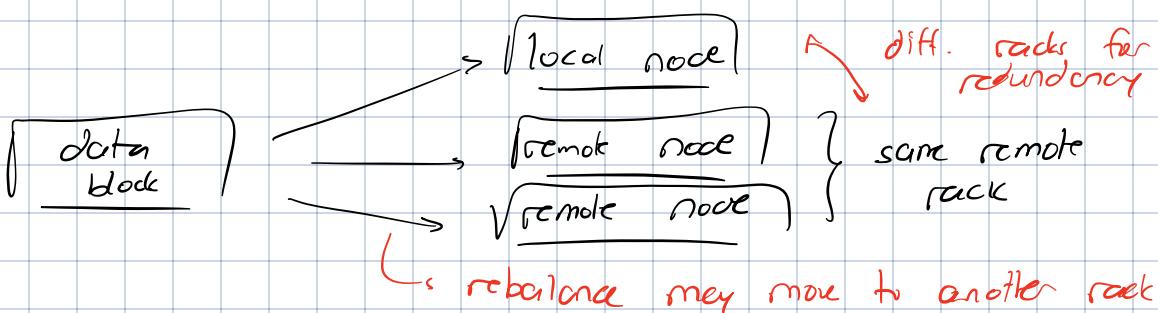
- ① Manage file system namespace
 - i) File \leftrightarrow blocks
 - ii) Block \leftrightarrow datanodes

} all in memory
- ② Cluster config. management
- ③ Block replication engine
- ④ File metadata & transaction log

Roles of datanode:

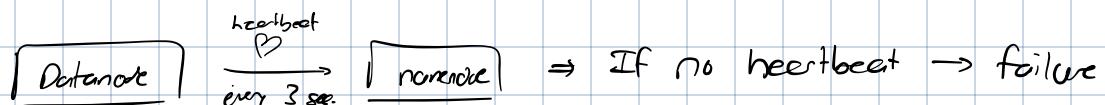
- ① Store data & metadata of blocks
- ② Report block status to namenode periodically
- ③ Forward data to other datanodes

Q1 Where exactly is data stored? Which datanodes? Default:



Client reads from closest block

Q2 How does namenode know if datanode failed?

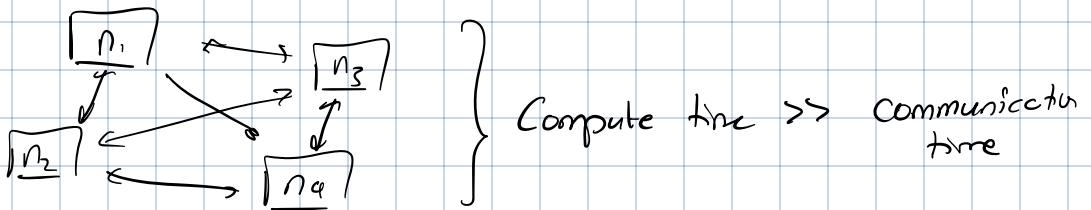


On fail, NameNode chooses DataNode for new replicas.

↳ Balances disk usage & comm. traffic

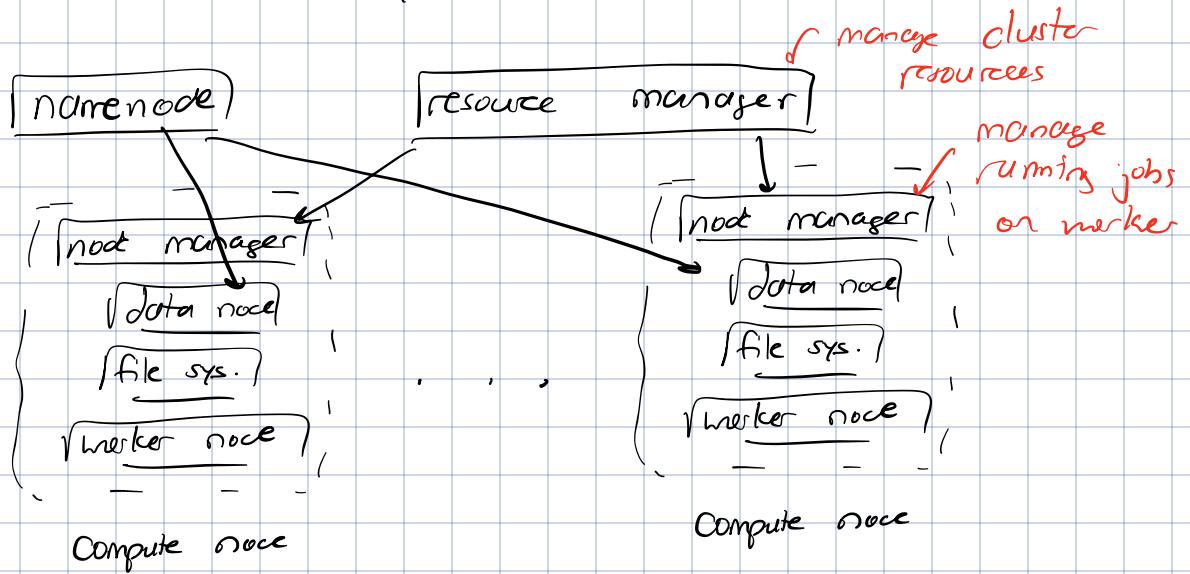
Q3 Which nodes actually do work on data?

For compute-intensive tasks, can distribute data to compute nodes



Problem for data-intensive comp: comm. time is too long!

Solution: do as much computation locally \rightarrow Hadoop assigns workers to node, that already have data



Algorithm Design

① Combiner design

B/c combiners may not run, we need to ensure that it matches the input & output types of mappe & reducer

Why do we want combiners? Combines are run locally to do file merge, reducing network traffic \rightarrow latency \downarrow

② Class variables

Sometimes, it can hold enough data in-memory. Be careful

③ In-mapper combine / in-memory combiner

Defn: perform reduce in a map task → network traffic ↓

Pros: speed

Cons: memory management required

MAPREDUCE To SPARK

Higher-level Programming

Problem w/ Hadoop: too much boilerplate & repetition

Companies created abstractions: FB Hive, Yahoo's Pig

Spark

Fast & efficient cluster computation engine compatible w/ Hadoop

Efficient

① Execution graphs & optimizations

② In-memory storage

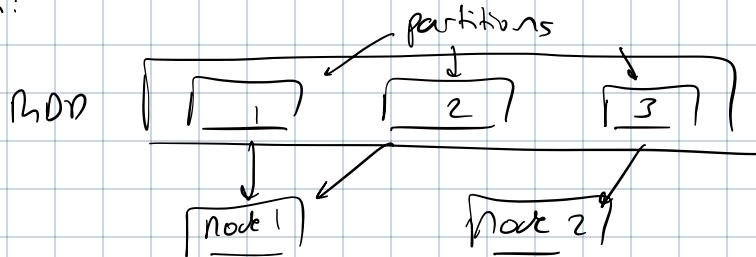
Useful b/c lots of good APIs

Spark performs ops. on RDD: resilient distributed datasets

◦ Defn: collections of objects spread across cluster

◦ Resilient b/c auto-rebuilds on failure as it tracks lineage data for RDD

◦ Diagram:



Programming w/ RDDs

① Creating RDDs:

sc.parallelize(...)

Spark context

OR sc.textFile("...")

⑦ Basic functions:

A: map:

$rdd.map(fnc)$ \Rightarrow run each elem through fnc.

B: filter:

$rdd.filter(pred.)$ \Rightarrow keep elements passing predicate

C: flatMap

$rdd.flatMap(fnc)$ \Rightarrow map each element \in RDD to 0/more other elem.

D: collect

$rdd.collect()$ \Rightarrow collects RDD into local obj.

E: take

$rdd.take(k)$ \Rightarrow emit 1st k elements

F: count

$rdd.count()$

G: reduce

$rdd.reduce(fnc)$ \Rightarrow merge elem. w associative func.

H: Save AsText File

$rdd.saveAsTextFile("hdfs://...")$

⑧ Key-value pair ops.

Spark's distributed reduce functions requires RDDs of KVPs

A: reduceByKey: $rdd.reduceByKey(-)$ \Rightarrow Each key has all assoc. vals reduce.

B: groupByKey: $rdd.groupByKey()$ \Rightarrow merge KVP into key: [values..]

C: sortByKey: $rdd.sortByKey()$

All ops. take a parallelism level operator

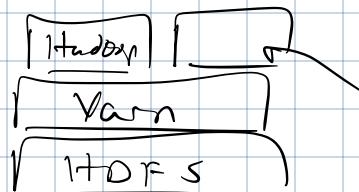
⑨ DAG:

Spark makes DAG for all ops. & uses it to perform efficient ops.

Data locality & partitioning is factored in

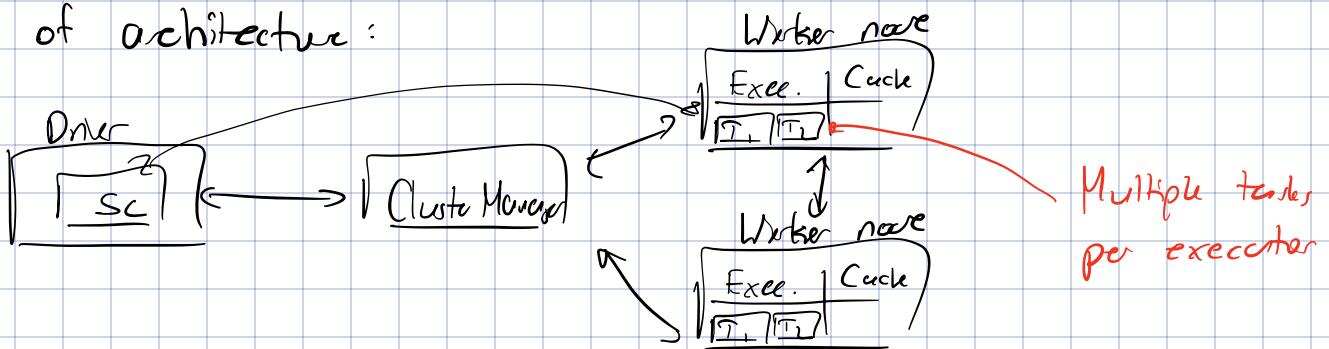
Spark Architecture

YARN is used to handle scheduling of jobs & resource management in Hadoop 2.0



Other distributed software,
like Spark

Diagram of architecture:



Driver sends Spark code on a per task level → memory problems.

Solution: broadcast

- Spark sends copy of value to executor, not just per task
- Broadcast values are read-only

To communicate back to driver, use sc. accumulator(..)

- Driver can inspect, workers can only write

Spark uses hash partitioner by default, but you can override.

Algorithm Design

① reduceByKey

- Reduces b/f & after shuffle!

② CombineByKey

- More fine-grained reduceByKey

- Signature:

$$\text{RDD. } [(k, v)] \cdot \text{combineByKey } (\text{create}, \text{append}, \text{merge}) \Rightarrow \text{RDD. } [(k, c)]$$

$\downarrow \quad \downarrow \quad \downarrow$
 $v \rightarrow c \quad c + v \quad c + c$

③ Aggregate By Key

- B/w Combine By Key & reduce By Key

$\text{RDD}[\{\langle K, V \rangle\}]$. aggregate By Key (init, append, merge) $\Rightarrow \text{RDD}[\{\langle K, C \rangle\}]$

Initial value
of type C \downarrow $C + V$ \downarrow $C + C$

④ Group By Key

Needed only if not reducing $\{\langle K, \text{List}(V) \rangle\}$

ANALYZING TEXT

Natural Language Processing

Aim: create a probabilistic model $P(w_1, w_2, \dots, w_k)$

- All LLMs are doing is given w_1, \dots, w_k , find prob. distribution of w_{k+1}
- Conditionally: $P(w_1, \dots, w_k) = P(w_1) \cdot P(w_2 | w_1) \cdot P(w_3 | w_1, w_2) \cdot \dots$

Using this to calculate $P(w_{k+1})$ is hard! Why? length [sentence] is unbounded so possibly LOTS of calculation

We can scale it down so we only use past n-words to construct prob.

- This is called n-grams
- Ex:// If bigram:

$$P(\text{I like ham}) = P(\text{I}) \cdot P(\text{like} | \text{I}) \cdot P(\text{ham} | \text{like})$$

Problem: unlikely words can make prob $\rightarrow 0$

- Solution:

① Robin Hood: every $P(n\text{-gram}) > 0$

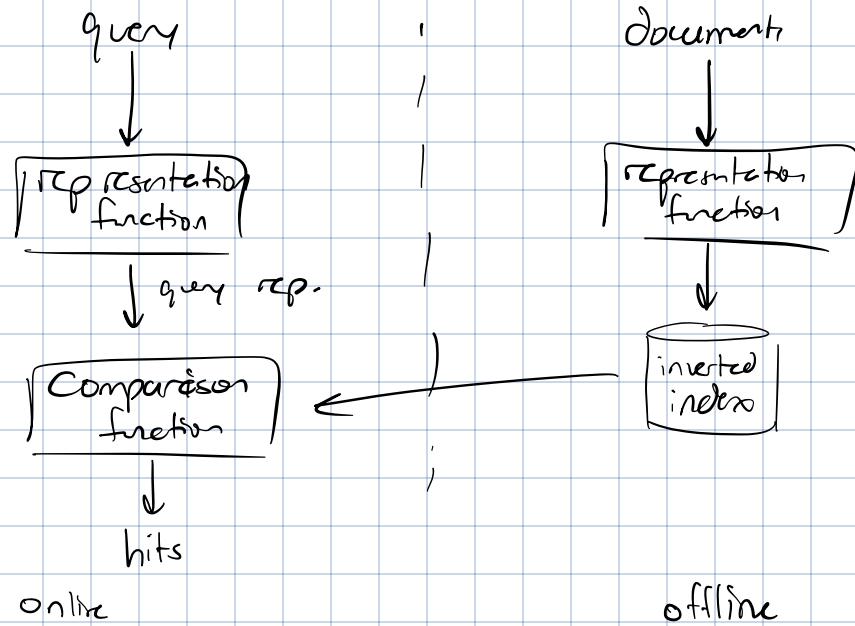
② Laplace smoothing: every count of a n-gram starts at 1

$$P(A, B) = \frac{\text{Count}(A, B)}{N} \rightarrow P_L(A, B) = \frac{\text{Count}(A, B) + 1}{N}$$

This requires adding N^2 to N (v being vocab size)

Searching & Indexing

Abstract information retrieval (IR) architecture



Representation function:

- ① Tokenize text, tokenize & case fold (lower case & Unicode → canonical form)
- ② Create embeddings

Inverted index: context → documents. (forward index: documents → context)

term1 → document1, document2, ...
 term2 → document2
 :
 } Postings list

Scaling:

° Heaps Law: $M = kT^b$ ↗ vocab size ↗ constant (0.4 - 0.6)
 ↗ # of docs. \Rightarrow vocab size grows unbounded
 ↗ constant (30 - 100)

° Zipf's Law: $f(k; S, N) = \frac{1/k^2}{\sum_{n=1}^N 1/n^2}$ ↗ rank ↗ exp. \Rightarrow Few elem occur v. freq.,
 Many elem occur v. infreq.

Map Reduce:

Mapper

def map(docid: Long, doctext: String):
 counts = Counter()

Reducer

def reduce(term: String, postings):
 p = list()

```

for term in tokenize(docText):
    counts.add(term)

for term, freq in counts:
    emit(term, (docid, freq))

```

Should
only be
for
big
freq
terms

for docid, freq in postings:
 \Rightarrow p.append((docid, freq))
 p.sort() ← Hadoop is good w/
 emit(term, p)

Delta encoding:

- Insight: if a term is common, delta of associated docs w/ each other is small
- Store only $\nabla(\text{doc}_i, \text{doc}_{i+1})$ for common terms
- V. useful if dealing w/ variable-length ints since it can be compressed
- Variable length ints: $\overbrace{1000\text{SLL}}^L$ bytes for x
 $\xrightarrow{\text{special}} 0$ if +, 1 if -

Elias-γ code: another compression format

- Assumptions: dealing w/ nat. # w/ no upper bound, smaller # more common
- Encoding a:
 - ① $N = \lfloor \log_2 x \rfloor$
 - ② Write N 0's
 - ③ Write x as $N+1$ -bit #

Decoding x

- ① Read 0s until 1, call it N
- ② Interpret $N+1$ bits as binary #

power law

Golomb Code: assume x is what we want to encode & decode, $x > 0$

1. $Z = \frac{\# \text{ of docs w/ term}}{\# \text{ of total docs}}$
2. $M = \lceil \frac{\log(2-Z)}{-\log(1-Z)} \rceil$
3. $Q = x-1/M \Rightarrow \text{quotient}$
4. $R = x-QM-1 \Rightarrow \text{remainder}$

- Each term gets own encoding scheme

Unary encoding: n is n 0s, then 1
 $\hookrightarrow S \Rightarrow 000001, 0 \Rightarrow 1$

Truncate binary:

- ① $k = \lfloor \log_2 n \rfloor$
- ② $U = 2^{k+1} - n$
- ③ First U code words:
 first v codes w/ length k
- ④ Last $n-v$ code words:
 last $n-v$ codes w/ length $k+1$

Ex:// Encode n=15 #'s in truncated binary

- ① $k = \lfloor \log_2 15 \rfloor = 4$
- ② $U = 2^{k+1} - 15 = 1$

③ First \cup codewords ($\{0, \dots, 14\}$) is cod \cup length k

$$\therefore v=1 \Rightarrow 0 \rightarrow \underbrace{000}_{k=3}$$

④ Last $n-k$ codewords: $k+1$ length

$$1 \Rightarrow 0010$$

$$2 \Rightarrow 0011$$

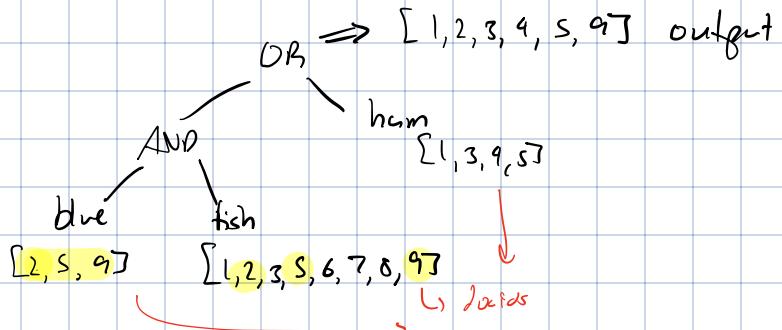
We can also collect term frequency for each document, but need to use a diff. approach

Indexing is a perfect problem for MapReduce

Searching & Retrieval

Boolean retrieval: take Boolean query \rightarrow query syntax tree \rightarrow for each leaf, look \in postings list to find docids \rightarrow apply query tree \rightarrow return hits (unsorted)

- Ex:// Q: (blue AND fish) OR ham:



Term-at-a-time: for each term, generate docs, then apply Boolean ops.

- We may need inverted sets for dealing w/ NOT operators

Document-at-a-time: run query tree against each document

- If index partitioned by document, this is nice \rightarrow each partition has all terms for a document, so no cross-partition queries needed

What if we want to sort output by relevance? Ranked retrieval

↳ No Boolean queries b/c we want to allow documents that might not have all terms

↳ Insight for relevance:

① If doc. has a term repeat many times (high term freq.), should be relevant

② Should be weighted according to commonality of term (doc. frequency)

TF - IDF:

total # of docs

$$w_{ij} = \text{tf}_{ij} \cdot \log N / df_i$$

weight of term i in doc j \rightarrow # of occurrences of term i in doc j
 ls # of occurrences of term i in doc j

$\log N / df_i$ \rightarrow # of docs containing term i

• Relevance of doc j: $\sum_{i=1}^{Nr} w_{ij}$

Document -at-a-time ranked retrieval:

- Algo:
 - for each doc:
 - score each query term \Rightarrow sum
 - accumulate best k hits via min-heap
- Con: can't terminate early, need to look at all docs.

Time: $O(n \log k)$
 Space: $O(k)$

Term -at-a-time ranked retrieval

- Algo:
 - collect hits & rankings for rarest term in accumulator
 - for each other term in query
 - if term \in doc:
 - remove doc from acc.
 - else:
 - add term ranking to overall ranking
- Early end possible, but uses memory

Tradeoff b/w 2. Commonly, docs partition by ID & sorted by quality \Rightarrow can do "early termination" document -at-a-time.

ANALYZING GRAPHS

Graph Review

Representing graphs:

① Adjacency matrix

$n \times n$ matrix ($n = \#$ of nodes). $M_{ij} = \begin{cases} 1 & \text{if node } i \text{ --- node } j \\ 0 & \text{o.v.} \end{cases}$

	1	2	3	4
1	0	1	0	1
2	1	0	1	1
3	1	0	0	0
4	1	0	1	0

Pros: matrix is great for math

Cons: wasted space w/ lots of 0s \rightarrow sparsity leads to wasted space

② Adjacency list

For i is list of out-neighbors of vertex i

$$1 : 2, 4$$

$$2 : 1, 3, 4$$

$$3 : 1$$

$$4 : 1, 3$$

Pros: smaller than adj. matrix & easy to find out-neighbors of v

Cons: difficult to find in-neighbors of vertex

③ Edge list

List all the edges

Pros: simple. Cons: not useful in finding neighbors

Web Data

Web is basically a huge graph, but quite sparse (most pages have 0/1 link)

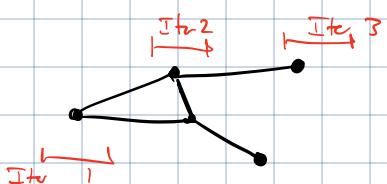
We use adj. lists for graph representation

Problem: find shortest path b/w vertex u & v :

- We could use Dijkstra's, but not parallel. Instead use parallel breadth-first search (pBFS).
- Algo:

$$\text{dist}(s) = 0$$

$$\text{dist}(v) = \begin{cases} 1 & \text{if } s \rightarrow v \\ 1 + \min_u (\text{dist}(u)) & \text{o.w. } \forall u \text{ s.t. } u \rightarrow v \end{cases}$$



- Hadoop

D.S.: { node n : (d: distance to n , adj list of n) }

Mapper:

```
emit (id, node)
for m in adjList:
    emit (m, node.d + 1)
```

Reducer:

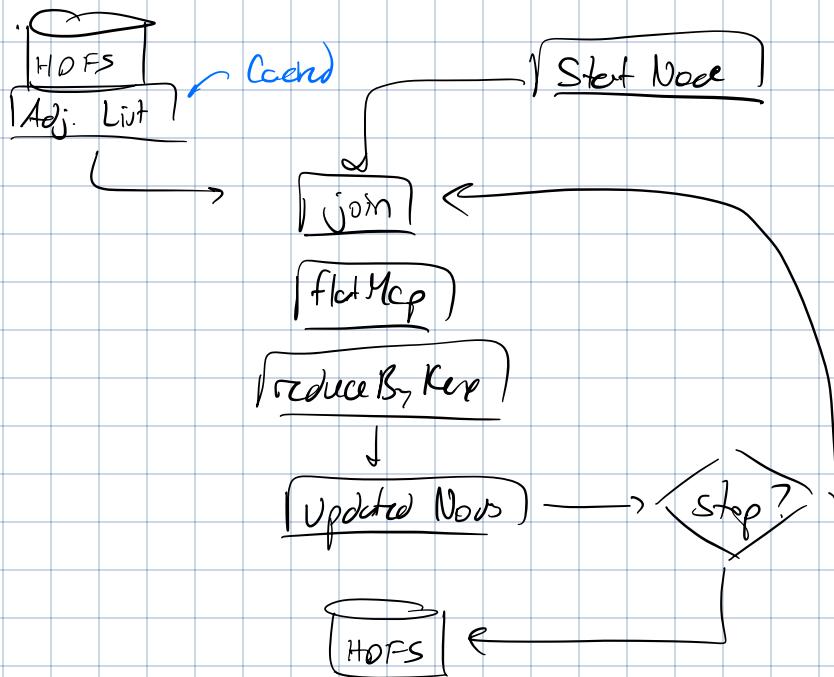
```
d = infinity
node = None
for o in value:
    if isNode(o): node = o
    else: d = min(d, o)
node.d = min(node.d, d)
emit (id, node)
```

- Can modify to actually give path:

- ① Message includes path that gives length d \rightarrow bigger message
- ② Node has a "prev" field \rightarrow requires reconstruction

- Problem w/ Hadoop: if we want to do this for all nodes, we need to read from disk a lot. Also need to keep sending graph struct to reducers every iteration

- Spark



Page Rank

Ranked retrieval fails b/c website w/ more terms \neq more relevant

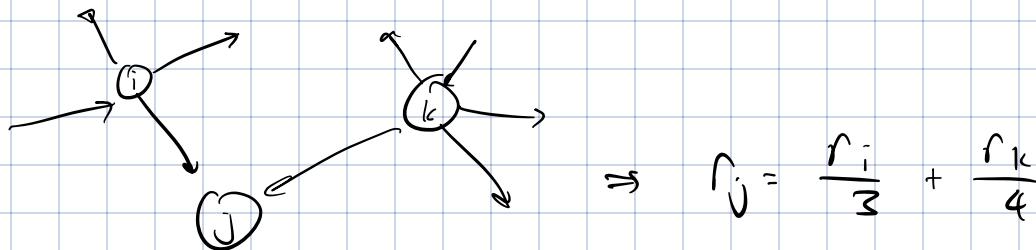
Instead, proxy will be # of in-links, w/ more important in-links weight more

If page j w/ rank r_j has n out-links, then each link conte

as r_j/n votes.

Then: rank of page $i = \sum$ ranks of in-links

Ex://

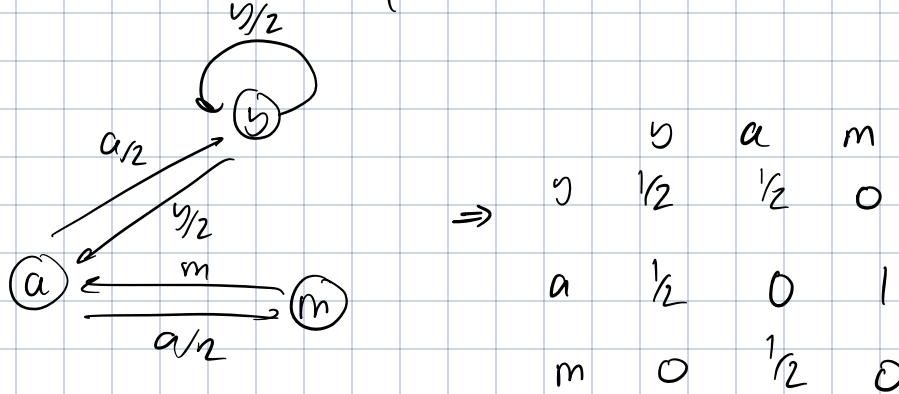


Stochastic adjacency matrix:

Let M be the matrix, page i has d_i out-links.

$$M_{ji} = \begin{cases} \frac{1}{d_i} & \text{if } i \rightarrow j \\ 0 & \text{o.w.} \end{cases} \Rightarrow \text{cols. sum to 1}$$

Ex://



Solving for ranks:

① Set $r_j = 1/N$ ($N = \# \text{ of nodes}$)

② $r'_j = \sum_{i \rightarrow j} r'_i / d_i$ until convergence

③ $r_j := r'_j$

Random walk interpretation: assume surfer is following links on web.

$\vec{p}(t)$: vector, i^{th} comp. is $P(\text{surfer on page } i \text{ at time } t)$

$$\therefore \vec{p}(t+1) = \underbrace{M}_{\text{stochastic matrix}} \cdot \vec{p}(t)$$

Important result: stationary distn. is unique & will be reached for graphs that satisfy certain conditions

Problems to solve:

① All out-links are w/in a group (spider-traps) \rightarrow importance is too high

\hookrightarrow Google soln: teleports

- Surfer follows random out-link w prob. of β , can also jump to random page at prob. $1-\beta$

Without this, ranks are too high for groups

② Some pages have no outlinks (dead-ends) \rightarrow importance "leaks out"

\hookrightarrow Soln: always teleport from deadends

Ex://

$$\begin{matrix} & y & a & m \\ y & \frac{1}{2} & \frac{1}{2} & 0 \\ a & \frac{1}{2} & 0 & 0 \\ m & 0 & \frac{1}{2} & 0 \end{matrix}$$

Dead-end

transform

$$\begin{matrix} & y & a & m \\ y & \frac{1}{2} & \frac{1}{2} & \frac{1}{3} \\ a & \frac{1}{2} & 0 & \frac{1}{3} \\ m & 0 & \frac{1}{2} & \frac{1}{3} \end{matrix}$$

Without teleports, columns are not stochastic

Google's Page Rank formula:

$$r_j = \sum_{i \rightarrow j} \beta \frac{r_i}{d_i} + (1-\beta) \frac{1}{N} \quad \Rightarrow \text{Assumes } M \text{ has no dead-ends}$$

↳ # of out-links for d_i

Page Rank w/ MapReduce

① Simple model: no random jumps & dead-ends

Map: node sends importance to out-links

Reduce: node sets importance to sum of received values

```
def map(id, n):
    emit (id, n)
    p = n.rank / len(n.adj)
    for m in n.adj:
        emit (m, p)
```

```
def reduce (id, msgs):
    n = None
    sum = 0
    for o in msgs:
        if o is node: n = o
        else: sum += o
    n.rank = s
    emit (id, n)
```

② Random jumps

```

def reduce(id, msgs):
    n = None
    sum = 0
    for o in msgs:
        if o is node: n = o
        else: sum += o
    n.rank =  $s * \beta + (1 - \beta) / N$  → Prob. of jumps to any random node
    emit (id, n)
  
```

③ Dead-ends

Option 1: add links to all nodes \Rightarrow too many msgs.

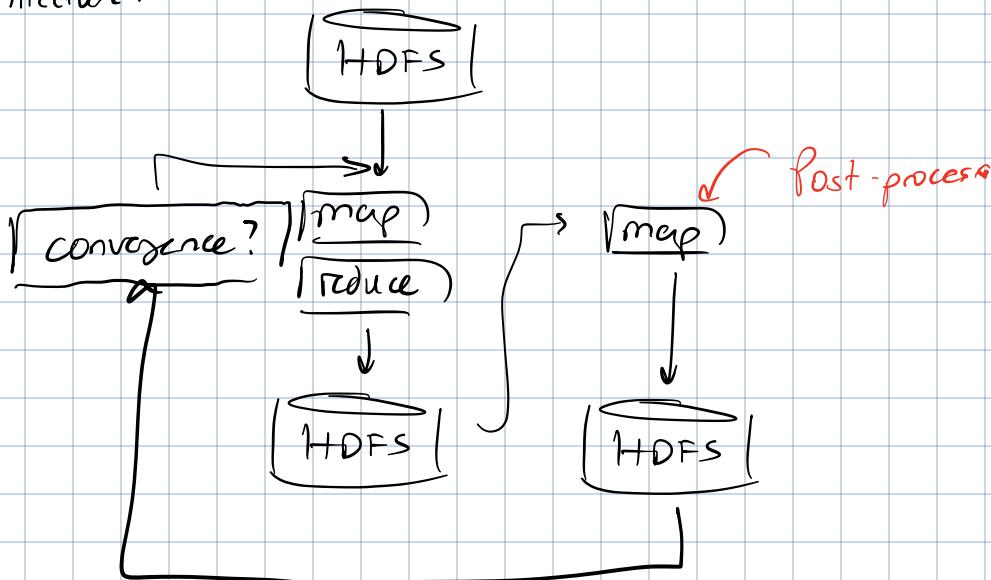
Option 2: post-process

R = sum of ranks of all nodes.

$1 - R$: deadends.

\therefore Add $\frac{1 - R}{N}$ to all nodes.

Architecture:



Alternative:

Mapper:

1. $\beta(\text{node.rank})$ to all out-links
2. $(1 - \beta)\text{node.rank}$ to "everyone"
3. node.rank to "everyone"

Reducer:

$$S \dagger = \text{"everyone"} \cdot \frac{1}{n}$$

Can we reduce the amount of space? Yup! Use log masses.

Q1: What is log mass?

$$\text{Let } a = \log m, b = \log n \Rightarrow m = e^a, n = e^b$$

$$\therefore m \times n = e^{a+b} \Rightarrow m \times n \sim a + b$$

Expansion:

$$m \times n = \begin{cases} b + \log(1 + e^{a-b}), & a < b \\ a + \log(1 + e^{b-a}), & a \geq b \end{cases}$$

Q2: Why is this helpful?

Ranks $\in (0, 1)$. $\log(\text{rank}) : (0, 1) \rightarrow (-\infty, 0)$

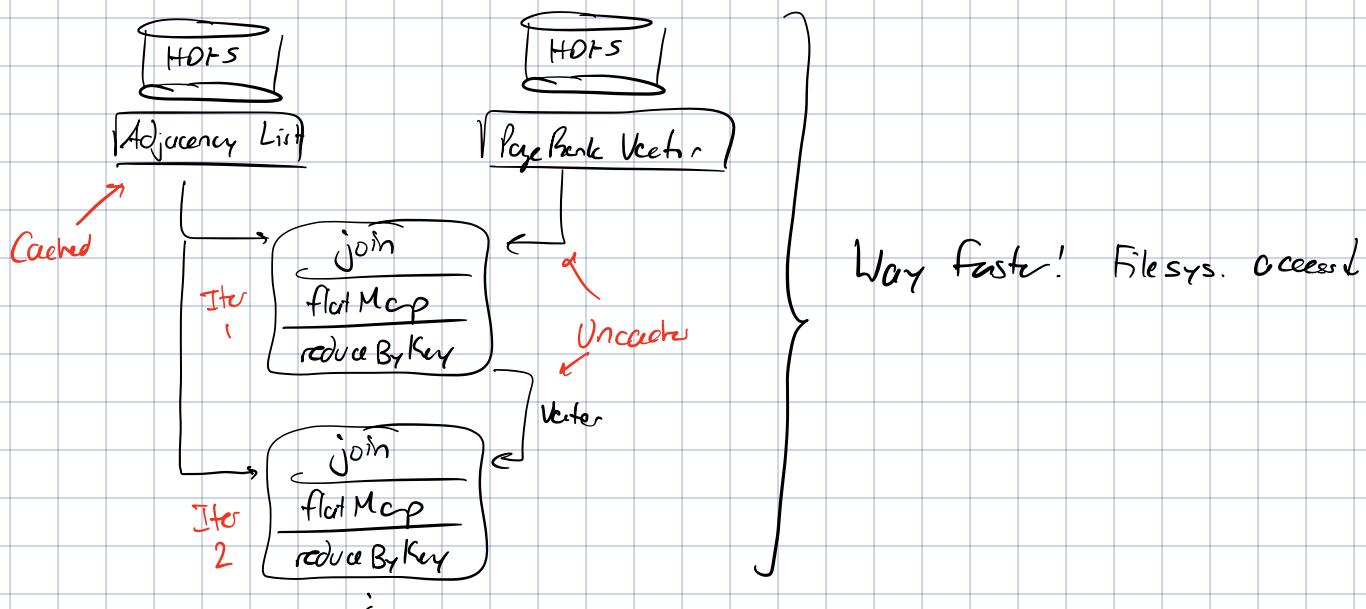
Uses more of the floating pt. range & won't have underflow b/c few vals will be close to 0

Page Rank w/ Spark

Problems w/ MapReduce:

- ① Sending adjacency lists w/ each iteration \rightarrow unnecessary messages
- ② Too much shuffling & file system access
- ③ Verbose
- ④ Each iteration of PageRank is a new job \rightarrow lots of time taken on startup

Spark architecture:



Page Rank Improvements

W/ Page Rank, we no longer have term spam that would elevate importance in trad. TF/DF algos.

Problem #1 w/ Page Rank: Spam Farming

- Post a website on a website w/ high rank \rightarrow your website rank \uparrow
- Alternatively, spider trap: a bunch of websites just point to 1 website
 - ↳ Random jump in Page Rank occurs accumulating too much rank in website, but still high
- Soln:
 - ① Add a "nofollow" attr. to links posted on popular websites to prevent rank inflation
 - ② Personalized Page Rank: can set trustworthy pages to start Page Rank

Q: How to determine if a website is trustworthy?

- Trust Rank
- 1. Set small set of ultra-trustworthy websites to source nodes in personalized pgc rank, set to 1
 - 2. Run personalized Page Rank \rightarrow every website has trust score $\in [0, 1]$
 - 3. Set trust threshold & cut off websites under threshold

NOTE: spam pages usually don't make cut b/c v. few trustworthy websites link to them

NOTE: personalized page ranks \rightarrow teleport to known good websites (source nodes) \rightarrow only good websites get rank

Need to ensure that starting set is trustworthy & big enough

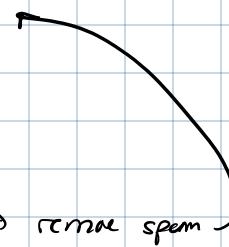
Algorithm to enhance trust set:

1. Run Page Rank

2. Select top k pages as seed

3. Run TrustRank

4. Select from websites w/ trust rank \geq threshold \Rightarrow remove spam



Q: How can we then determine definitively if a page is spam?

r_p := PageRank of page p

$r_p^+ := \text{PageRank of } p \text{ from trusted sources}$, random jumps only lead to trusted source

$r_p^- := r_p - r_p^+ \Rightarrow \text{Contrib. of low trust pages to rank}$

$s_p := \frac{r_p^-}{r_p} \Rightarrow \text{Spaminess}$

$s_p \uparrow \rightarrow \text{more likely } p \text{ is spam}$

DATA MINING / MACHINE LEARNING

Supervised ML

Defn: give model training set where you already label data w/ soln. Model learns & then applies to new data

Data is represented as vectors of features.

↳ Features can be sparse: lots of 0s

} Also called embeddings

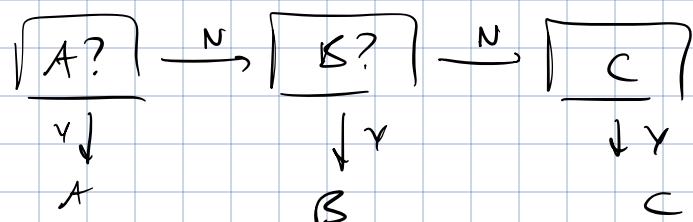
You want to choose good features! Don't choose everything

Process:

- ① Acquire data: usually want as much as you can for better perf.
- ② Determine labels of data: can crowdsource, bootstrap or exploit user behaviour
- ③ Choose model.

We will use binary classifier as building block.

More complex models use cascade of binary clf. to do multi-label classifier



④ Train model

Want to find function f s.t. it minimizes

$$\frac{1}{n} \sum_{i=1}^n l(f(x_i, \theta), y_i)$$

↳ l is a loss func.

We often do this computationally via gradient descent

↳ For each component x of θ , find $\frac{\partial l}{\partial x}$ & go downhill \rightarrow compute θ_{i+1}

$$\therefore \theta' = \theta - \gamma \nabla L(\theta)$$

↳ step ↳ gradient

How do we avoid getting stuck in a local minima?

1. Momentum

2. Stochastic gradient descent: randomly select subset of training data

In binary CLF, we use logistic function:

$$f(x; w, b) : \mathbb{R}^d \rightarrow [0, 1] = \sigma(w \cdot x)$$

\downarrow

$$\sigma(z) = \frac{e^z}{e^z + 1}$$

Then:

$$l(y, t) = \frac{1}{2} (y - t)^2$$

This bounds the loss & it is differentiable

Q: Can we run gradient descent on MapReduce?

Yes, but we need to run multiple jobs (1 per iteration). This incurs:

1. Startup time
2. Need to retain state b/w iter.
3. Bottlenecked jobs (1 reducer)

Q: Could Spark do better?

Yup! It covers data, no shuffles & data only does 1 reduce.

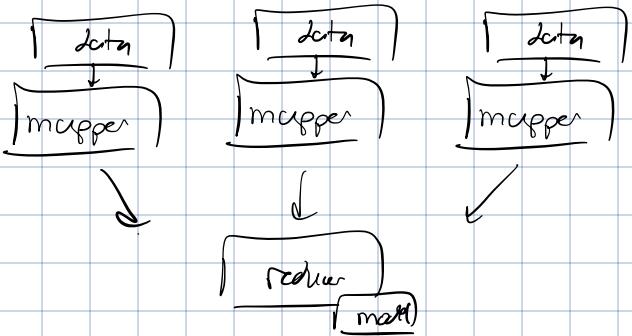
Ensemble learning: use multiple models indep. & combine predictions

↳ Works b/c errors are uncorr., $P(\text{all models wrong}) \downarrow$, variance reduced

Can use this for stochastic gradient descent!



Qn: mapper is prior & reducer is learner:

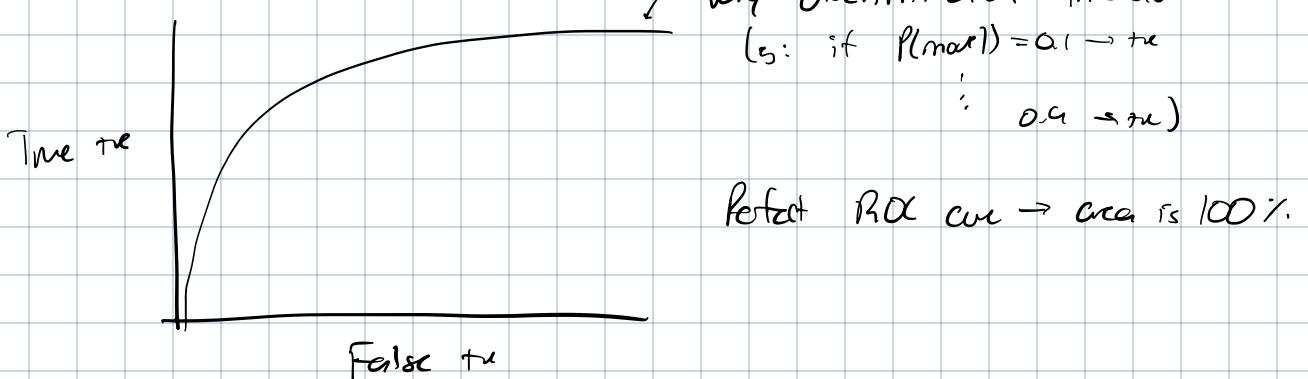


Model Evaluation

Recall that loss measures fit of model; we want to see how well model does on unseen data

Many metrics:

① ROC curve:



Can use k-fold cross-validation to continuously test & train model

Embeddings

Problem w/ feature embeddings is that they can be really big & sparse

Solution: reduce dimensions

↳ One soln. is to hash \rightarrow modulo by some prime.

Nearset neighbors problem

Obj.: given high-dim. data points x_1, \dots, x_n and distance function $d(x_i, x_j)$, find all pairs (x_i, x_j) s.t. $d(x_i, x_j) < s$

① Naive approach:

Compute $d(\dots)$ for all pairs $\rightarrow O(n^2)$

② Locality sensitive hashing approach

Locality sensitive hashing: items that are "close" will have hashes that are "close"

Also:

1. Make LSH table

2. Use buckets of width w . Buckets can overlap

3. Most values x_i, x_j s.t. $d(x_i, x_j) < s$ will have at least 1 bucket in common

How to make locality sensitive hash algorithm (specifically for sets)

↳ Jaccard similarity:

Let C_1 be set of n -grams of doc 1, C_2 for doc 2

$$\text{sim}(C_1, C_2) = \frac{|C_1 \cap C_2|}{|C_1 \cup C_2|}$$

$$d(C_1, C_2) = 1 - \text{sim}(C_1, C_2)$$

For multisets (elem. can appear multiple times), union & intersect defn. changes:

$$\text{If } C = A \cup B, \text{ then } C[k] = A[k] + B[k]^{\substack{\leftarrow \\ \# \text{ of times } k \text{ in } A}}$$

$$\text{If } C = A \cap B, \text{ then } C[k] = \min(A[k], B[k])$$

For text, we will use n -grams as embeddings. size of n depends on doc size

Recall that sets & vectors can be represented as each other:

Sets \rightarrow vectors?

1. Label all elem. of \mathcal{U} set from $1 \rightarrow n$

2. If set contains i , then vector at $i = 1$, o.w. 0

Min-hashing: let C be a set of n -grams and π is an enumeration func.

$$h(\pi; C) = \min_{i \in C} \pi(i)$$

Property: $P(h(C_1) = h(C_2)) = \text{sim}(C_1, C_2)$.

↳ Proof: $y = h(C_1 \cup C_2) \Rightarrow y = h(C_1)$ or $h(C_2)$.

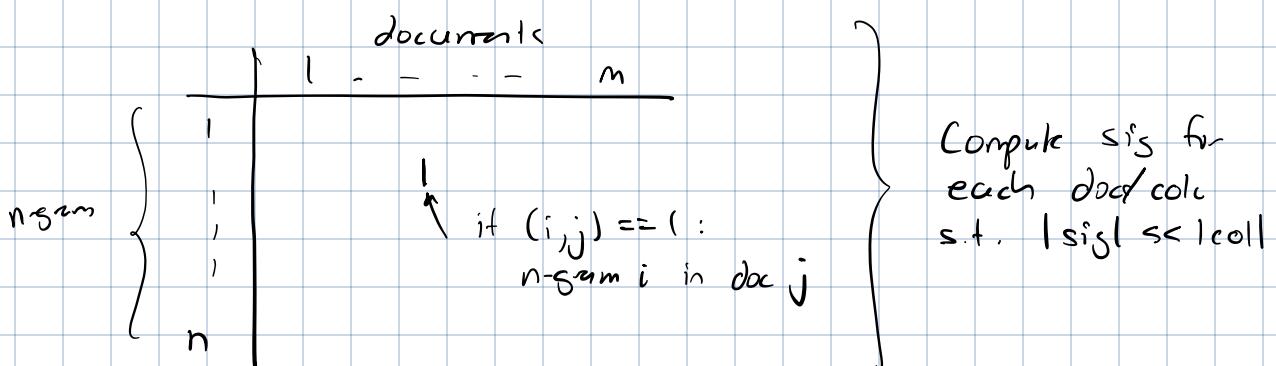
This is because C_1 or C_2 might be smallest value of π , we don't know since it depends on π

Only same if in intersection of C_1 & C_2

$$\therefore P(h(C_1) = h(C_2)) = \frac{|C_1 \cap C_2|}{|C_1 \cup C_2|}$$

This property indicates that $h(\pi; C)$ is an LSH

We want to do the following:



Define: $\text{sig}(C)[i] = h(\pi_i; C) \Rightarrow$ min-hash for i^{th} perm out of K.

↳ $|\text{sig}| = K \times 4 \text{ bytes} \ll 1\text{coll}$

Ex://

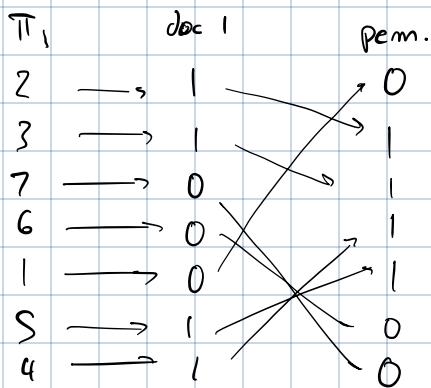
	π_1	π_2	π_3
2	4	3	
3	2	4	
7	1	7	
6	3	2	
:	:	:	
:	!	!	

	input matrix			
	doc 1	doc 2	doc 3	doc 4
ns 1	1	0	1	0
ns 2	1	0	0	1
ns 3	0	1	0	1
ns 4	0	1	0	1
:	0	1	0	1
:	0	1	0	0

Compute signature for documents.

Steps:

1. start w/ doc 1. Reorder accrdg to π_1 ,



Problem: calculating π functions is too much!

2. Find row index w/ first 1:

This row, index 2.

3. Repeat for π_2 & π_3 , ...

4. Combine to form $\text{sig}(C_1)$

Define $\text{sim}(C_1, C_2) = \% \text{ of entries that are equal}$.

↳ For signatures, just compare $\frac{\# \text{ of indices in same order}}{\text{sig. size}}$

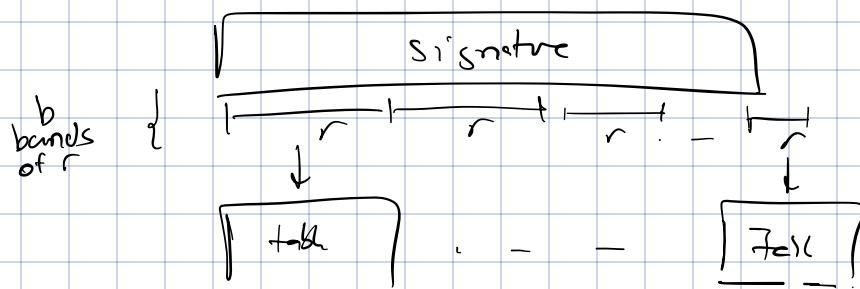
We can create a k-universal hash function:

$$h(x) = c_1 + c_2 x + \dots + c_k x^{k-1} \pmod p$$

↳ Ensures that $h(x_1), \dots, h(x_n)$ do not correlate.

Another problem: signatures are vectors \rightarrow still need to compare if signatures of 2 docs are similar, which requires n^2 comparisons (if $n = \text{length of sign.}$)

Soln:



If 2 signatures collide in any 1D table, signatures are candidate pair

Ex// Assume 100k docs, each w/ signature of 100 bits. Choose $b = 20$ bands and $r = 5$ integers/band. Want to find doc pairs that are $s=0.8$ similar

Assume $\text{sim}(C_1, C_2) = 0.8$

$$\hookrightarrow P(C_1 \text{ & } C_2 \text{ identical in 1 bin}) = (0.8)^5 = 0.328$$

$$P(C_1 \text{ & } C_2 \text{ not sim in all 20 bins}) = (1 - 0.328)^{20} = 0.00035$$

This tells us that 1/3000 of our similiar doc pairs are false neg.
So we can find 99.965% of sim. documents.

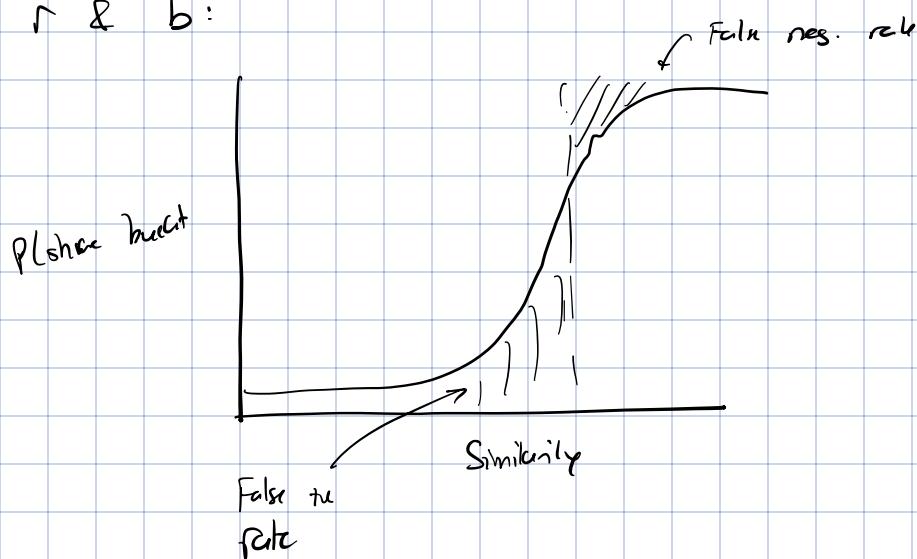
What if $\text{sim}(C_1, C_2) = 0.3$

$$\hookrightarrow P(C_1 \text{ & } C_2 \text{ identical in 1 bin}) = (0.3)^5 = 0.00243$$

$$\hookrightarrow P(C_1 \text{ & } C_2 \text{ not identical in 1/20 bins}) = 1 - (1 - 0.00243)^{20} = 0.094$$

\hookrightarrow ~5% of docs of 0.3 sim are considered candidate pairs

To pick r & b :



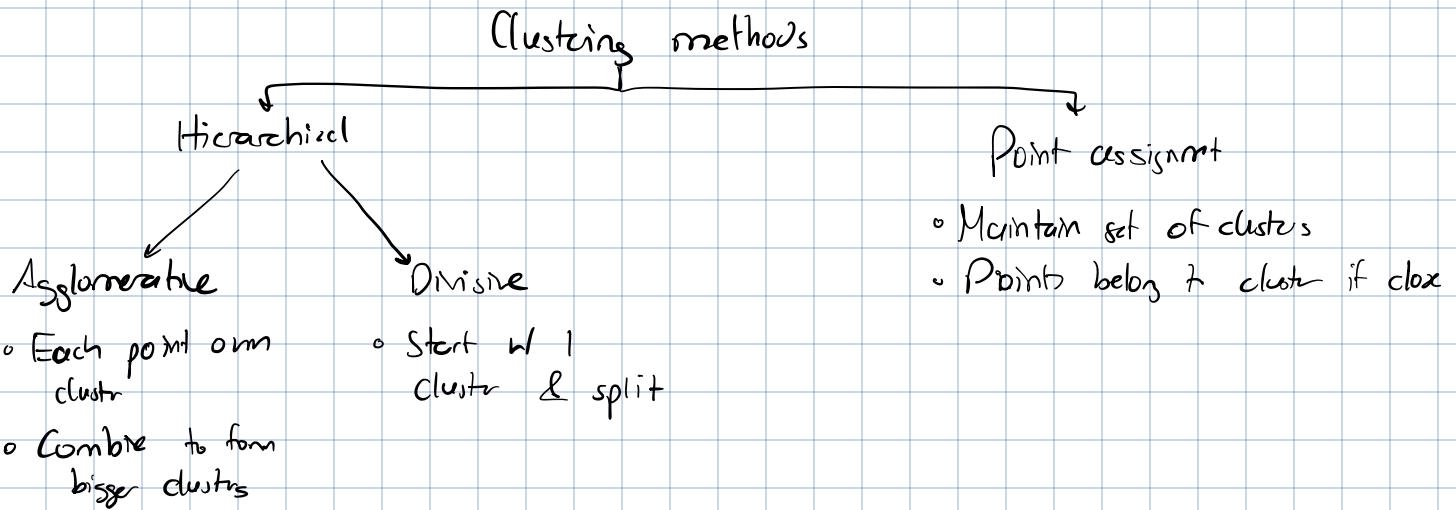
Summary:

- (1) Doc \rightarrow n-grams \rightarrow integers.
- (2) Generate universal hash functions
- (3) Compute signature vector for each doc
- (4) Turn R, B for false +ve & false -ve rate.
- (5) Hash each of B bins for each doc to find pairs

Clustering

Problem: given set of points, group into clusters

↳ Difficult b/c points often defined in high-dimensional spaces



① Hierarchical clustering

1) How to represent a cluster of points?

In Euclidean case, cluster has a centroid (avg. of data points)

If not Euclidean, choose centroid: point closest to other points

2) How to measure fitness of clusters.

Distance b/w centroids for Euclidean case.

For non-Euclidean:

1) Intercluster distance: min (distance b/w any 2 points from each cluster)

2) Cohesion

A: Diameter of one cluster

B: Avg. distance b/w points in cluster

C: Density of cluster

② K-means clustering

Assume Euclidean space & pick k to be # of clusters

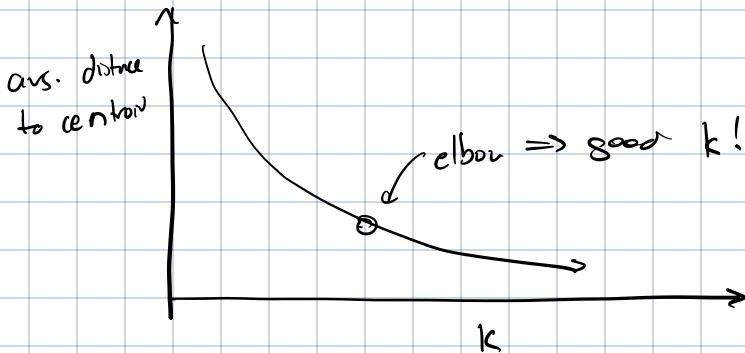
Initialize k clusters by picking k points.

Algo:

1. For each point, place in cluster whose centroid it is closest to
2. Update centroid locations for all clusters.
3. Reassign points

repeat
until
convergence

How to select k?



Map Reduce implementation

def setup :

clusters = loadClusters()

def map (id, vector):

emit (clusters.findNearest(vector), vector)

def reduce (clusterId, values): Iterable[Vector]

for (vector ← values):

sum += vector

cnt += 1

emit (clusterId, sum / cnt)

average distance
per cluster

Spark implementation

1. Pick random centroids

2. Create RDD w/ (cluster #, data point) → partition by key & create

3. Reduce by key & map → new centroids

4. Create new RDD by reassigning centroids.

Track # of points that changed via accumulator

Faster b/c

1. All data in mem.
2. Points that don't change cluster they don't go to another worker

If accumulate > 0.

RELATIONAL DATA

DB workloads

OLTP

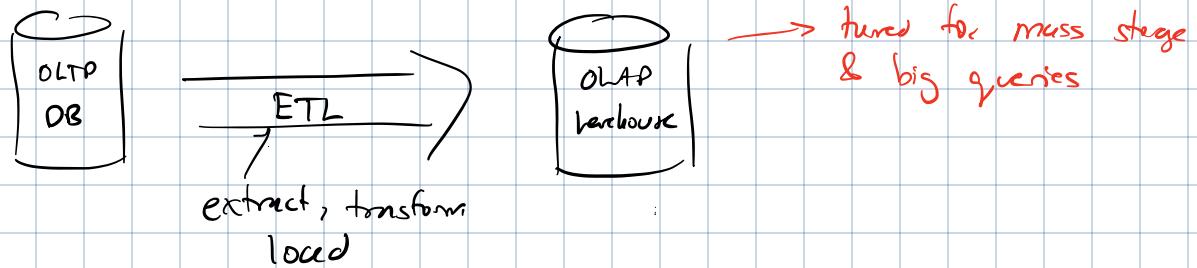
- Most apps use OLTP
- User facing: fast, concurrent
- Usually small set of common queries
- Random reads & small writes.

OLAP

- This is for BI & data mining
- Batch workloads, low concurrency
- Complex, ad-hoc analyses.
- Full table scans, big data

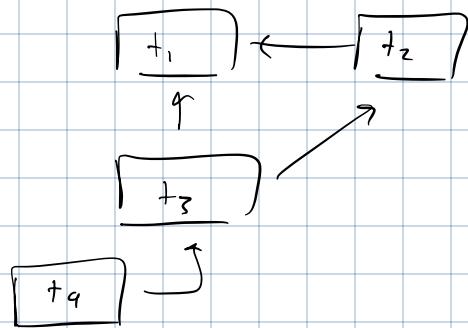
| DB doing both causes a lot of problems

Soln: Data warehouse



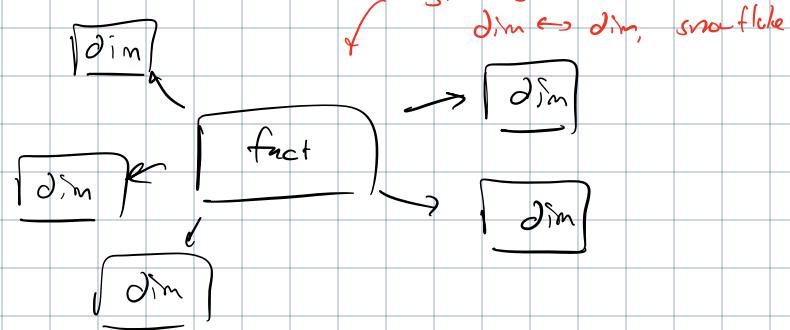
Schemas will be diff:

OLTP:



Captures relations b/w objects.

OLAP:



Fact tables: have unique info per row

Dim tables: foreign keys to fact table
that is not unique to a single row.

Ex:// Sale of widget → fact table (color w/ price, item ID, ...)

Store sold, manufacturer, ... ⇒ dim table

At large scale, OLAP warehouse can just be Hadoop. Analysts use Hive for SQL → Hadoop

Data lake: schema-less dump of data w/ mix of raw & curated data

Data lake $\xrightarrow{\text{ETL}}$ warehouse

Q: How does a query get executed on big data warehouses?



RDBMS will do query \rightarrow AST \rightarrow execute (no Java int. repr.)

Relational Algebra

① Selection (σ)

Equivalent to SQL: `SELECT ... WHERE ...`

MapReduce implementation: map-side filter

\hookrightarrow b/c it is a map operation, it can be pipelined!

\hookrightarrow Bottleneck will be I/O of loading from text file

② Projection (π)

Equivalent to SQL: `SELECT col1, ..., coln ...`

MapReduce impl.: map-side filter

\hookrightarrow Need to remember column mappings after selection (e.g. $col4 \rightarrow col1$)

③ Rename (ρ)

SQL: `select col1 as ...`

Not important in MapReduce b/c column names don't matter

④ Union (\cup)

SQL: `(select ...) union (select ...)`

MapReduce: use MultipleInputFile class \rightarrow 1 Mapper class / input file

⑤ Difference (-)

SQL: `(select ...) minus (select ...)`

MapReduce:

1. Multiple Input File class

2. Each mapper emits: $((data, mapperId), \dots)$
3. If $S_1 - S_2$, sort S_2 tuples b/t equal S_1 tuples
4. Reducer: Remember last S_2 tuple & emit only S_1 tuples not equalling S_2 tuple

→ (a₁, S₁) → emitted by reducer
 (a₂, S₁) → emitted by reducer
 (a₃, S₂) → remembered by reducer
 (a₃, S₁) } → not emitted
 (a₃, S₁) }

⑥ Cartesian product (\times)

SQL: select table1..., table2... from table1, table2

MapReduce: cross joins are extremely expensive → lots of data

- If it is unavoidable, Hadoop will cross file partitions & custom reader will emit tuples to mapper

⑦ Aggregation:

SQL: COUNT, MIN, MAX, SUM . . .

MapReduce: do it on reduce si

Relational joins

3 types of joins: inner (or 1-to-1), one-to-many, many-to-many

Joining on MapReduce:

① Hash Join (also known as broadcast, side-loaded, replicated)

Assumption: t₁ is small enough to fit on single node of memory

Process:

1. Every mapper loads t₁ as hash table
 2. Each mapper joins its partition of t₂ against t₁
 3. Map emit
- } Fast!

② Map side join

If files sorted by key & split in some partitions, can join row by row rather than crossing entire files

Exactly like cross join, but you should have file sort & partition cond.

③ Reduce-side join (shuffle / repartition join)

Mappers send join key + data → reducers get tuples in join key sort order so reduce joins

Join types:

A: 1-1 joins \Rightarrow 1 key per t_i , so easy to join

B: 1-many joins \Rightarrow hold 1 tuple in. reduce memory from t_1 (ensure it comes 1st though key ordering) & then join w/ t_2 tuples

↳ Key: (join key, origin) (origin ord $\Rightarrow t_1 < t_2$)

C: Many-many joins: how to hold all tuples in memory from $t_1 \rightarrow$ memory problems

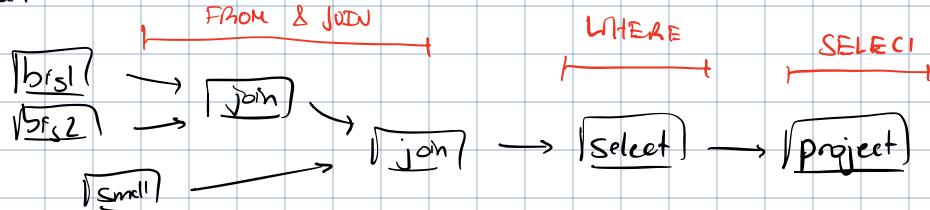
SQl to Hadoop

Know how to do this!

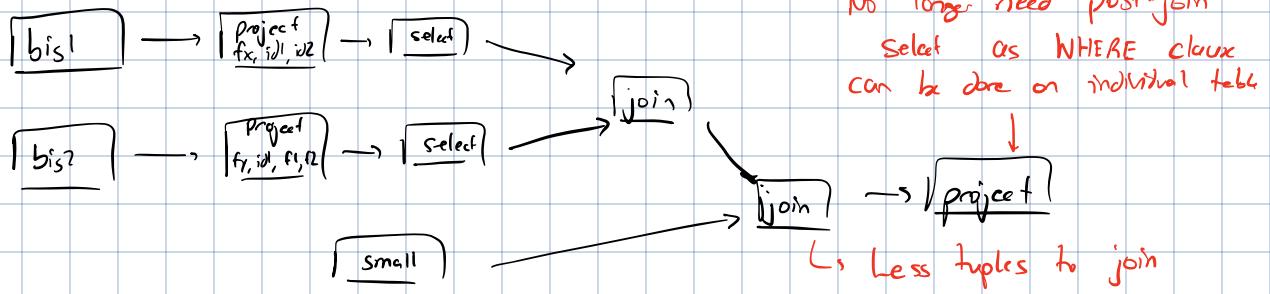
Ex:// Convert SQL to Hadoop:

```
select big1.fx, big2.fy, small.fz  
from big1  
join big2 on big1.id1 = big2.id1  
join small on big1.id2 = small.id1  
where big1.fx = 2015 and  
big2.f1 < 40 and  
big2.f2 > 2
```

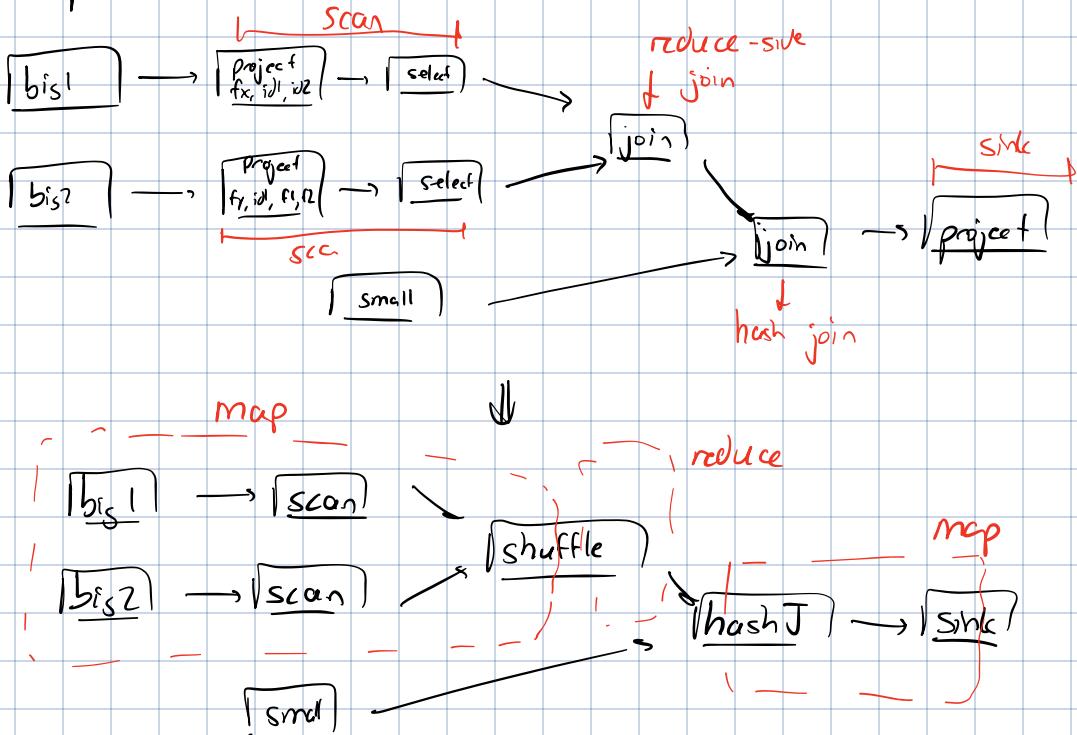
① Logical plan



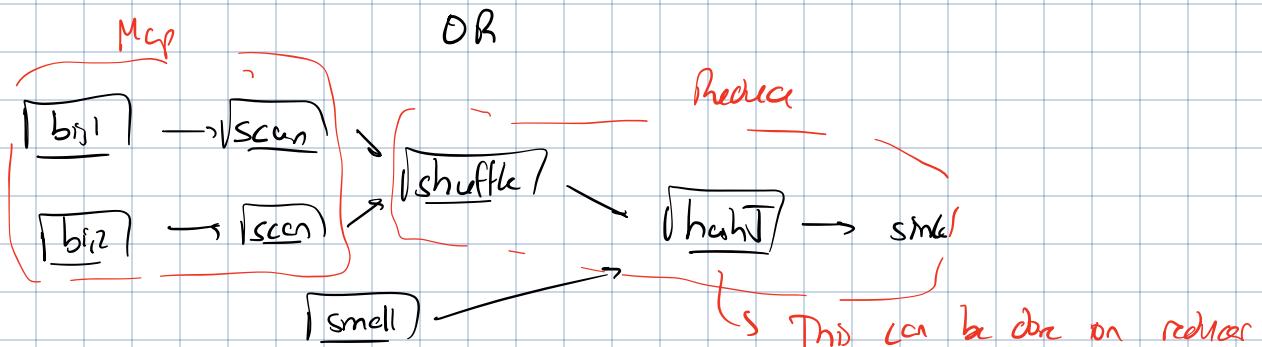
② Optimize logical plan



③ Physical plan



2 Map Reduce jobs



1 MapReduce job

We know table schemas from metadata

SQL → Spark

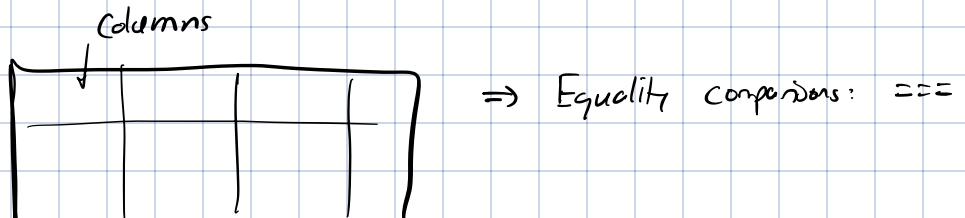
Most relational operators are already RDD transformations

Join:

- Spark will automatically choose map-side or reduce-side depending if it's a narrow dependency (no shuffling) or wide dependency (shuffling needed)
- If hash join needed, use a broadcast var

Spark also has Spark SQL which will transform queries into Spark operations

- Uses DataFrames as base data struct instead of RDD.



- RDD → DF:
`rdd.toDF("col1Name", "col2Name", ...)` Can also pass in column schema
- To query:
`df.createOrReplaceTempView("tableName")`
`spark.sql("sql query using tableName")`

Row vs. Column Storage

Reading in text is 'slow'. Convert to binary format & include a schema to convert binary-coded columns to logical columns.

Methods of storing data

Row store

row1, row2, row3
+

[col1 val] - - [coln val]

- Easier to modify record
- Requires reading a lot of data for processing

col values
usually repeated!

Col store

[col1 doc1 / col1 doc2] . . . [col2 doc1] . . . []

- Only read necessary data
- Underlie more complicated
- Can compress data
- Vector execution & compile query work well

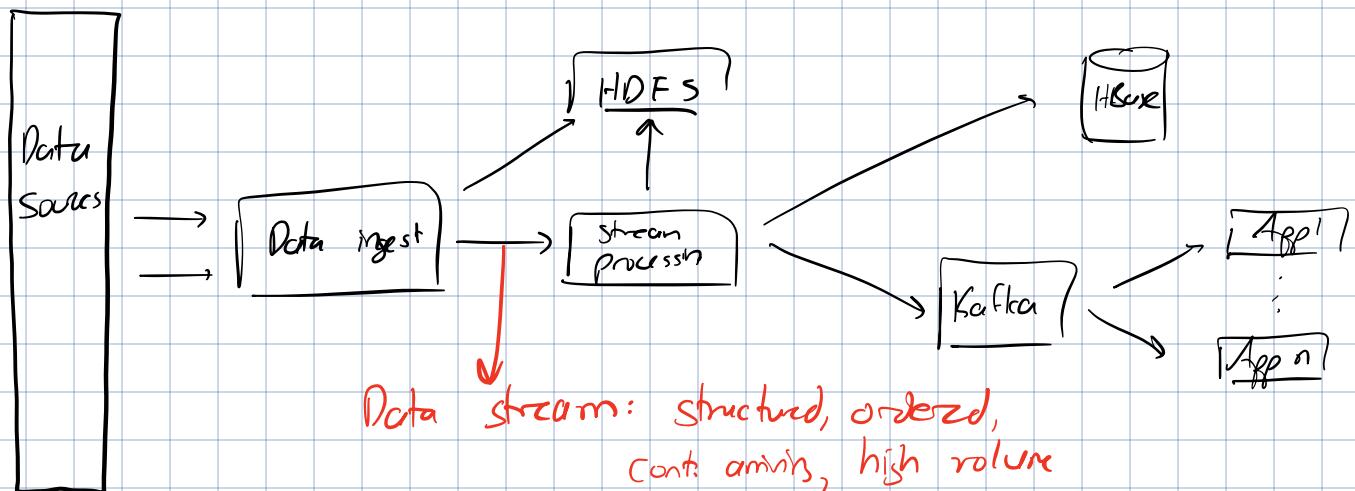
How can we use binary input in Hadoop? Parquet!

- Input: Parquet Input Format
- Output: Parquet Output Format

STREAMING

Stream processing: process data as it arrives, like real-time

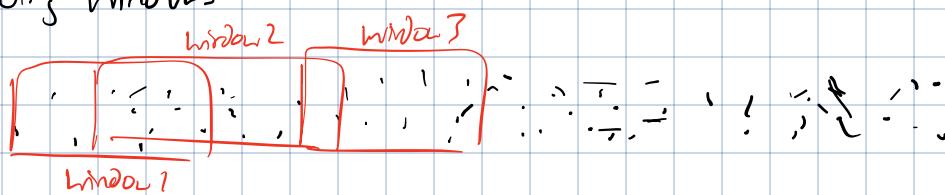
How does it fit in the modern data arch:



Q: How do you do reduce-ops on continuous data?

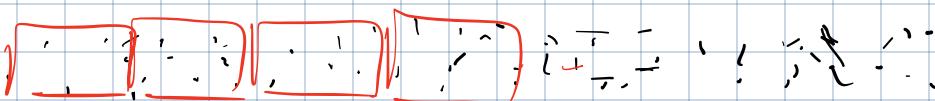
A: Define windows (10-minute, 1-hour, ...)

① Sliding windows



A window of x minutes updates y seconds would update y seconds to most current timestamp & have x minutes of data before y

② Tumbling windows



No data in common b/w windows

Windows can be based on:

1) Ordering attribute: time

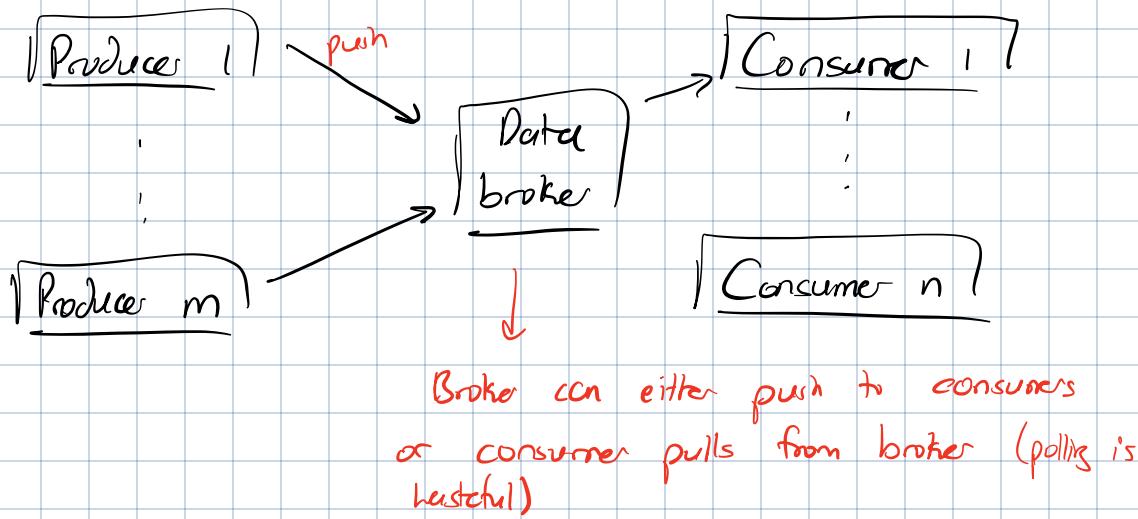
2) Counts: window has n data points b/f new window

3) Punctuation: ingest apps define window end, but its window size becomes unpredictable

Stream processing challenges:

- Latency requirements
- Memory storage
- Burstiness
- Load balancing
- Data consistency
 - ↳ Exactly once
 - ↳ At least once (dups possible)

Stream processing often works on pub-sub:



Spark Streaming

Tumbling window of extremely small duration → turn into RDD

Needs ssc (Spark streaming context) to perform opps on RDDs (seq. known as DStream)

Can process 6 GB/sec, higher throughput than Storm

Data Brokers

Kafka is very widely used for brokering. MQTT is for IoT uses

Kafka is fast (2M writes/sec). Why?

- ① Writes almost always go to OS caches
- ② Reads are fast b/c fast transfer from cache to network socket

Kafka topics are partitioned & replicated for fault tolerance

↳ Messages are hashed & balanced to diff. partitions

We want idempotence, so UVIDs attached to each message. Broker rejects messages if UVID seen b/f.

By keying a message, you guarantee that all msgs w/ same key get to same partition, which is ordered.

Multiple consumers can read a single partition & each consumer remembers how much it read

↳ Consumers can be grouped to read from single topic, but each consumer reads 1 partition

High Volume Streams

① Sample data randomly: select S values randomly

1) Reservoir sampling: start 1st S values, then start k^{th} w/ prob. $\frac{S}{k}$,
discard already stored

↳ Theoretical result: $k \geq S$. $\forall i \in [1, k]$, $P(V_i \text{ in reservoir}) = \frac{S}{k}$

Multiset Ops

Q: How do we know how many items in set

A: HyperLogLog counter

Observation: hash(item) \rightarrow vector of 32 bits

↳ 1/2 of items have code w/ 0., 1/4 w/ 00...

We can record longest string of prefix 0s & estimate 2^x

Q: How do we know if an item is in a set

A: Bloom filter

Create a vector of m bits & k unique hashing func.

put function: hash x & flip bits of k hashes in m-vector

contains function: \cap & check if all of K hashes is 1

No false negatives, but sometimes false positives (can be tuned)

Q: How do I count frequency of x in set

A: Count-min sketch

Like Bloom filter, but each hash func. has own m-vector bit vec.

Increment each vector hash elem

To get # of times x seen, take minimum of vector elems of hash

$$\min \left\{ \begin{array}{c} A[h_1(z)] \\ \vdots \\ A[h_q(z)] \end{array} \right\} = \dots$$

MUTABLE STATE

Problems w/ RDBMS:

- ① Hard to scale out \leftarrow relations b/w tables
- ② Doesn't support semi-structured

Alternative: NoSQL

↳ Features:

- ① Horizontal scaling
- ② Replicated & distributed data
- ③ Not ACID (eg. eventual consistency)
- ④ Schemas flexible

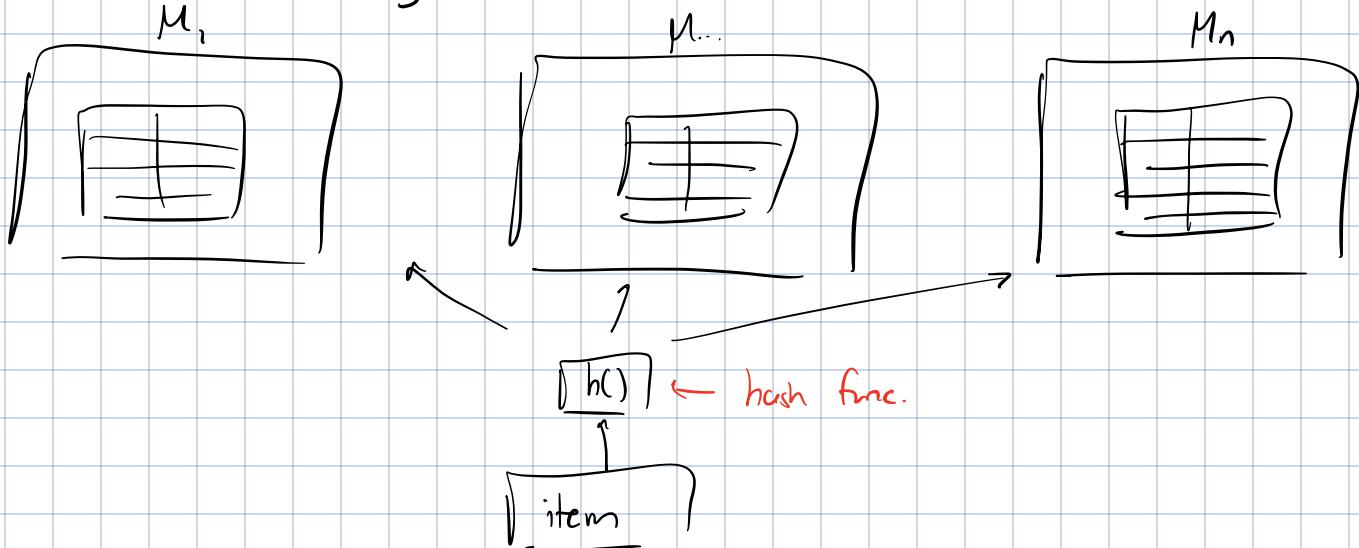
Follows BASE:

- ① Basically available: should be avail. most times
- ② Soft state: state may change even w/out input
- ③ Eventually consistent: updates may cause inconsistency, but will work out

Types of DBs: KV stores, column stores, document-based, graph

Key Value Stores

Goal: distributed hashing



Consistent hashing:

- ① Assign each node a GUID from hash func. h
- ② For key k , compute $h(k)$ ← same hash func.
- ③ successor(k): first $\text{GUID} \geq h(k)$

Why is this approach good? Each node has at most $\frac{(1+\epsilon) \cdot k}{n}$ keys.
 ↳ If node goes down, only removes $O(k/n)$ keys

How to make successor(k) fast? Chord distributed protocol

- ① Node sdn: link nodes in doubly linked circle & check all nodes until cond. matched
 $O(n)$ worst case
- ② Finger table: each node holds predecessor, successor & finger table:

successor	node	Finger table
$h+2$	1	
$h+9$	1	
$h+8$	1	
:		

Also:

- 1) Find $h(k) \rightarrow$ go to node

2) In node, look at finger table & so far as you can what exceeds $h(k)$

3) Repeat until successor found.

$O(\log n)$

To handle new nodes, we have to do following:

1) Create new finger table for new node

2) Update other node finger tables: background process

Assumption: peer-to-peer nodes, most systems are not like this.

Memcached: distributor KV store, no replication of data

↳ Acts as a cache (no persistence)

↳ Every node in cluster acts independently.

↳ hash func: $h(k) = k \% \text{NUM_SERVRS}$

Redis: can handle more than just KV

↳ Distributed: nodes can communicate \rightarrow replication \rightarrow redundancy

↳ Persistent

↳ Hashing: $\text{bucket}(k) = h(k) \% 16384 \rightarrow$ cluster coordinate which bucket belongs to which node.

Buckets can be added/reassigned to diff. nodes based on load / node health

Document Stores

Idea: key $\longleftrightarrow \{ \dots \}$ (JSON-like doc.)

Many systems use multi-value concurrency control for eventual consistency

key : $\boxed{t_1: \text{doc1}} \rightarrow \boxed{t_2: \text{doc2}} \rightarrow \dots$

- Since it is a linked list, reads & writes can be concurrent!
- Read: look at latest committed version
- Write: append to linked list

Graph DBs



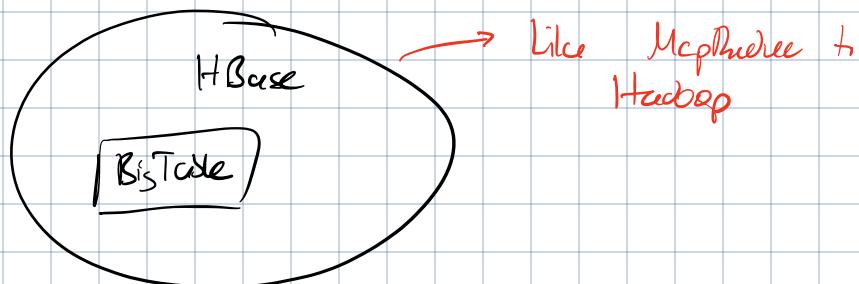
Used if you have a lot of relations (e.g. social media)

Common tech: Neo4j

Column Stores

Google BigTable: (row key, column key, timestamp) \leftrightarrow <bytes stream>

\hookrightarrow Bloom filters to check if rkey, ckey exist

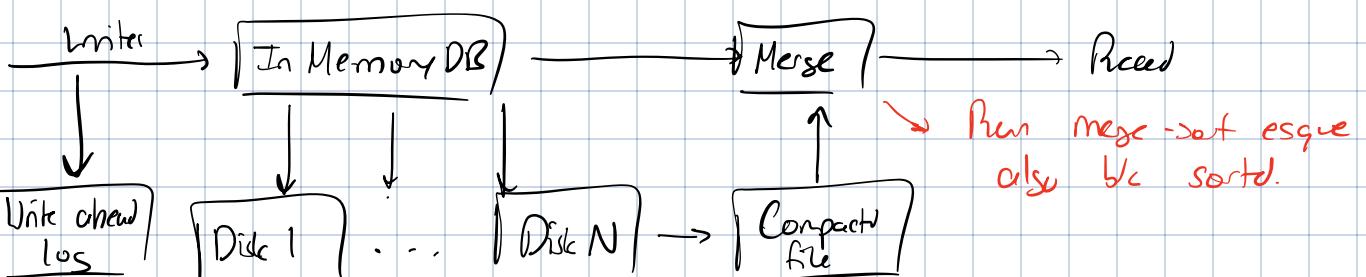


Backed by HDFS

Rows are maintained in sorted lexicographical order \rightarrow efficient row scans

Columns are grouped into families (Ckey: Cfamily + Cqualifier)
↑
in family

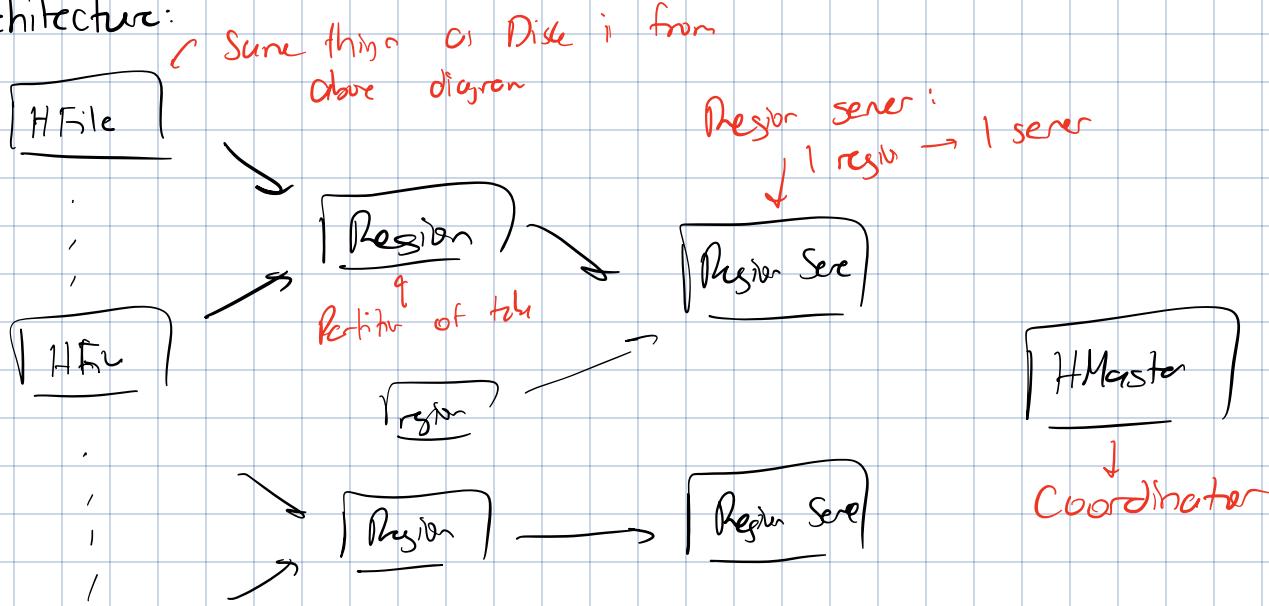
Underlying data structure: log structured merge tree (LSMT)



1) Replay if data loss

2) Can sync to other replicas
Why not copy to 1 disk file?
Appendix will break sorted order

Architecture:



Consistency

Consistency: all nodes return same value

Availability: function w/out all nodes being available

Partition tolerance: || network being split

CAP theorem: pick 2 from above cannot have all 3

↳ We want dist. sys. to work during partition tolerance. Pick b/w cons. or avail.

↳ SQL prefers consistency. NoSQL prefers availability

Types of consistency:

① Strong: after update, all subsequent access returns updated value

o 2-phase commit protocol

1) Coordinator makes sure everyone can commit

2) || tells subordinates to commit data

3) || ensures data got committed

} can be roll backed

Blocking & slow

- Distributed consensus protocol (e.g. Paxos)

Nodes can be proposers (advocate for data change), acceptors (decide if change should go) and/or learners (remember values)

Paxos algo:

- 1) Proposer selects proposal # $n >$ all prev. proposals
- 2) Proposes n to quorum (\geq majority) of nodes
- 3) Nodes either:
 - a) Promise to reject all future proposals $< n$
 - b) Reject proposal
- 4) If quorum sends promises, then proposer tells everyone to accept

② Weak: after update, some will return new value, others return old

◦ Eventual: after some time, all nodes return updated value

If no partitions, having tradeoff b/w consistency & latency (high consistency requires time)

HBase always keeps consistency no matter what

- Transactions only allowed to update 1 row
- Since 1 region \rightarrow 1 region server, it will be consistent