

big data.
definition.

possible solution { scale up / vertically expensive
scale down / horizontally

failure.

develop program
model

cloud computing = services through internet.

3 services: infra. platform. software

computing services development env applications
CPU, memory, network

2 deploy models { public
private.

5 characteristics

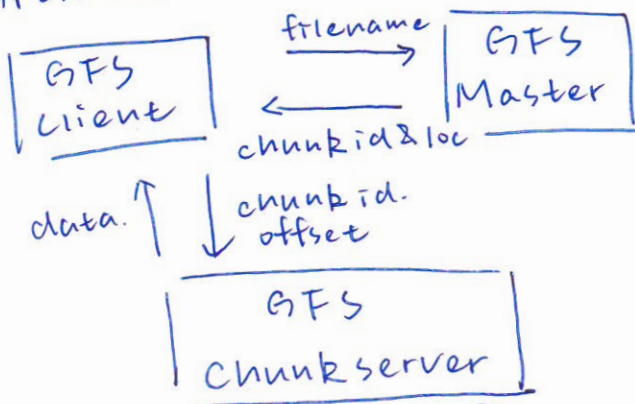
Lec 2 storage

file system on single machine: X scale.
large files

↓
Distributed file system

GFS.

Architecture.



Master: ① metadata (x3) ② heartbeat msg
Chunkserver.
Client

chunk size: 64 mb.

{ ~~adv~~ pro = ↓ metadata
↓ intereaction
con = wasted space ↑

Consistency Model

{ consistent
defined.

operations.

{ read
update
write
append
delete

Lec 3 storage. NoSQL.

Database revolution.

No DBMS \rightarrow hierarchical 1st \rightarrow Relational DB ACID 2nd \rightarrow NoSQL NewSQL 3rd \rightarrow BAg

SQL problems: scalability.

ACID properties: Atomicity, consistency, Isolation, Durability

Theory consistency \rightarrow strong consistency
eventual consistency

CAP theorem: consistency, availability, partition tolerance

BigTable HBase (CP). Cassandra (AP)

NoSQL BASE properties: Basic Availability, soft state. Eventually consistent

NoSQL Data Model \rightarrow key-value column-oriented. e.g. BigTable, HBase, Cassandra.
document graph. e.g. MongoDB

Big Table Data Model

Table columns: column family + qualifier



rows } tablet

} tablet

> Table

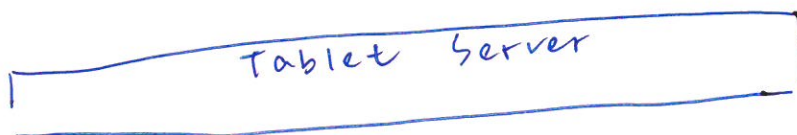


timestamp

Arch.

Master

client



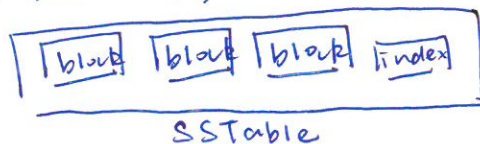
Master: ① assign tablet ② balance load ③ monitor schema changes

Tablet server: ① manage tablet

Building blocks: GFS + chubby (lock service)



} tablet



SSTable

Lec 5 Parallel Processing.

Problem: x load data in single computer
challenge: ① split computation ② manage failure.
MapReduce.

Programming model.

{ map phase: (key, value) pair \rightarrow intermediate () pair
shuffle phase: group key value pairs by key
reduce phase: grouped pairs \rightarrow final result.

execution:

files split \rightarrow master & workers \rightarrow master assign work
map worker output inter pairs \rightarrow buffered in memory
 \rightarrow write to disk.

master forward loc \rightarrow reducer worker

\rightarrow notify master \rightarrow wake up user program

Fault tolerance: heartbeat

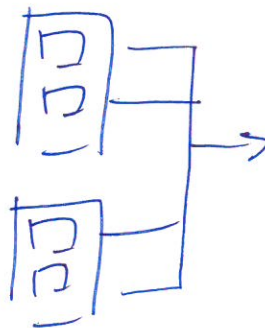
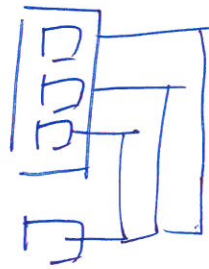
map { in progress: \rightarrow re-execute.
completed: }

reduce { in progress: re.
completed: result stored.

Algorithm
join

map side join.

reduce side join.



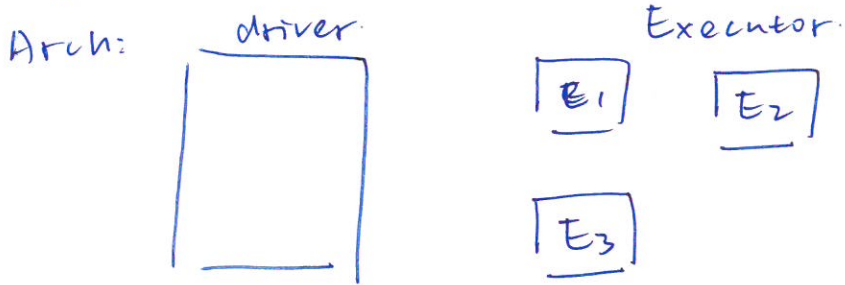
Problem: a chain MAPReduce.

\downarrow

Flume Java

Lec 6 Spark 1
motivation: data sharing (MR in slow storage)

↓
Spark: keep things in memory



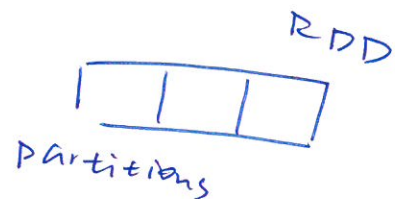
Driver: ① interact with user. ② distribute work with executor.

Executor: ① execute ② report state

Programming Model:

RDD (Resilient distributed dataset)

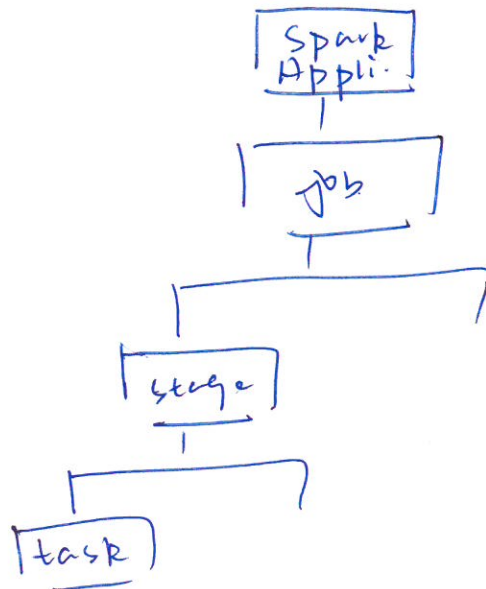
operations — transformations
 \ actions



Lineage ○ → ○ → ○ logical execution plan
 ↓
 ○

dependencies (edges) — narrow
 \ wide

Fault tolerance: re-compute



LeC7 Structured Processing = Spark SQL

MR: problem: no schema. X query

↓
Hive: + structured data on MR.

exe: parsing query to tree → generating logical plan

→ physical plan

↓
Shark: ~~modify~~ run on Hive.

limitation: X integrated with Spark. X optimized

↓
Spark SQL.

RDD: distributed data collections.

X know schema.

↓
DataFrame: structured data.

↓
X type checked by compiler.
(X compile time type safety)

Dataset:

base RDD + DF.

compile type safety + optimizer.

Execution: code → logical plan → physical plan

optimization — catalyst = reorder.op.
prune partition

↓ Tungsten: ↑ CPU memory.

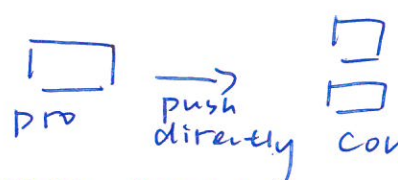
lot of heap memory,

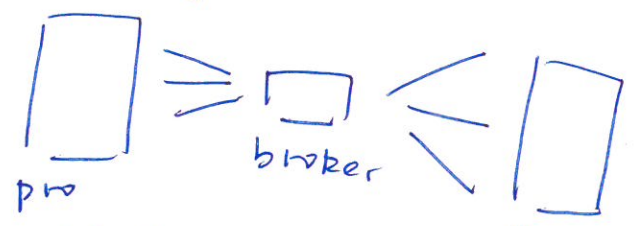
Lec 8. Streaming ~~processing~~ ^{storage} = log.

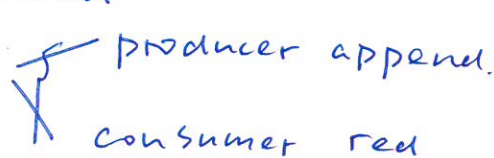
streaming data: unbounded input

storage problem:

message system: spread info from producer to consumer.

direct messaging: 
drawback: consumer crash
msg too fast.

message broker: 
typical: delete msg once consumer consumed

log-based: store in log 
producer append.
consumer read sequentially


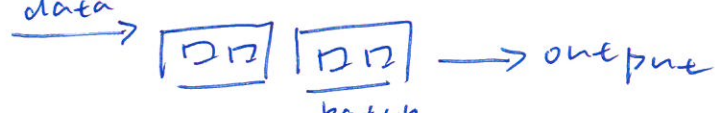
partitioned: hosted on diff machines
seq num within each partition

topic: partitions with same type msg

Kafka: replicate partition → fault tolerance

zookeeper: manage broker & consumer.
keep offset

streaming processing

1. continuous processing: 
micro-batch: 

2. Record-at-a-time API: low-level full control

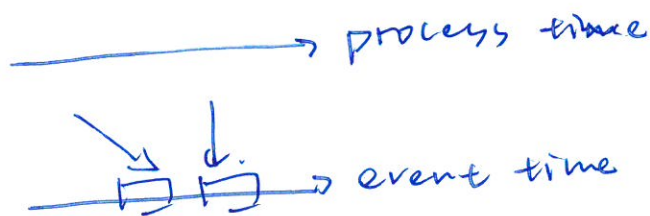
Declarative API: what to do v.

3. Event time: actually occurred

Process time: received.

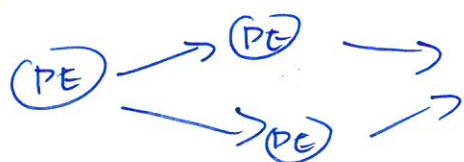
4. windowing — tumbling: all ~~delete~~
sliding: oldest delete

windowing — by processed time: fixed time window
 — by event: handle out-of-order events



Processing Model

Data flow graph =



PE (processing element)

Tasks — stateless: x require prior knowledge
 easy parallelized

stateful: involve diff tuples

Logical plan:

Physical

vertex = PE	edge = streaming connection
vertex = process	transport connection

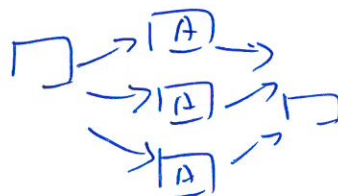
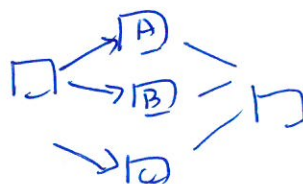
MAP logical → physical

1. Parallelization — pipelined

task



data.



2. Fault Tolerance:

delivery — at-least-once
 — exactly-once

recovery — active backup: (processing node + backup node)

same input

passive backup: checkpoint

upstream backup: x backup. store until Ack

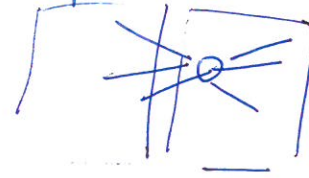
Lec 10 = Streaming processing = spark streaming

DStream (Discretized Stream Processing)

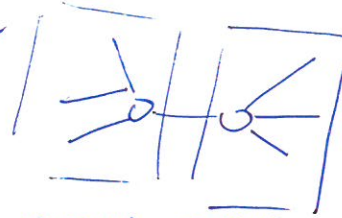
Lec 11 Graph processing.

Challenges: partition data, & computation

Graph partition — edge cut



vertex-cut



Graph processing = process large graph over

① Think like a vertex / Graph parallel

② compute own ③ depend on others

Pregel: synchronous

Execution: supersteps/iteration.

read (s-1) send (s+1)

Fault tolerance: workers revert to checkpoints

limitation: determined by slowest

GraphLab = asynchronous. shared memory (lock)

Fault tolerance: periodical halt + synchronize

PowerGraph = gather-apply-scatter
syn/asyn.

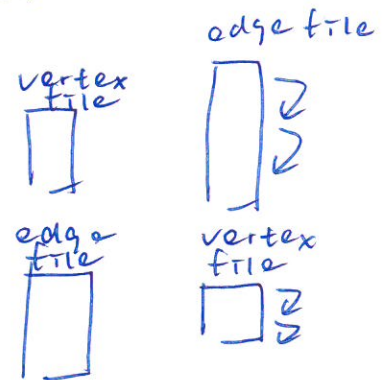
② Think like an edge

vertex-centric breadth first

edge-centric

problem: random access on vertex

streaming partition graph



③ Think like a table

motivation: restrict types

GraphX: Graph parallel in Spark
+ data parallel

Lec 13 Resource Management

CPU + RAM.

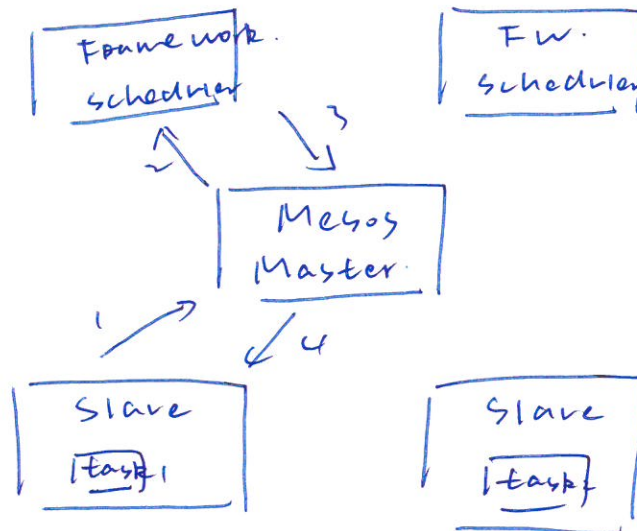
Mesos.

↳ Fine grained sharing: tasks level

coarse grained sharing: jobs/framework level

✓ Resource offer.

Arch:



Allocating resources

single resource — fair sharing. $\frac{1}{n}$

min-max.

weighted min-max

multiple resource — asset fairness (weights).

dominant resource fairness

YARN