

Welcome!

CS 744: SNOWFLAKE

Shivaram Venkataraman

Fall 2020

Prepare

→ Reading lecture

Prepare notes

→ Try to solve previous year exams

ADMINISTRIVIA

Open ended

→ Pros / Cons for various designs

→ Compare various systems

- Assignment 1 grades out! → Canvas. Contact Saurabh

- Assignment 2 by mid-week →

- Midterm this week! → Thursday. The exam will be posted on Canvas

- Project Proposal Peer review

↓
Assigning one
other project
for you to review

BBCollaborate → for clarifications

PDF upload to Canvas → Practice this!!

You can also type it out

AEFIS FEEDBACK

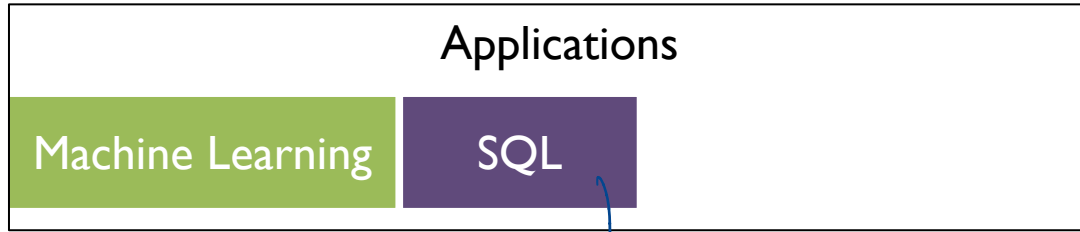
4 / 78 responses!

How has your experience been reading papers?

Are the lectures useful for learning?

How are the discussion groups? Did you get to know students in the class?
Would it help to have the same group each time?

Anything else we could improve for the second half?



→ SCOPE : Relational API

Analytics → Language
Integration

↓
DAG of operators
run on execution
engines like Spark

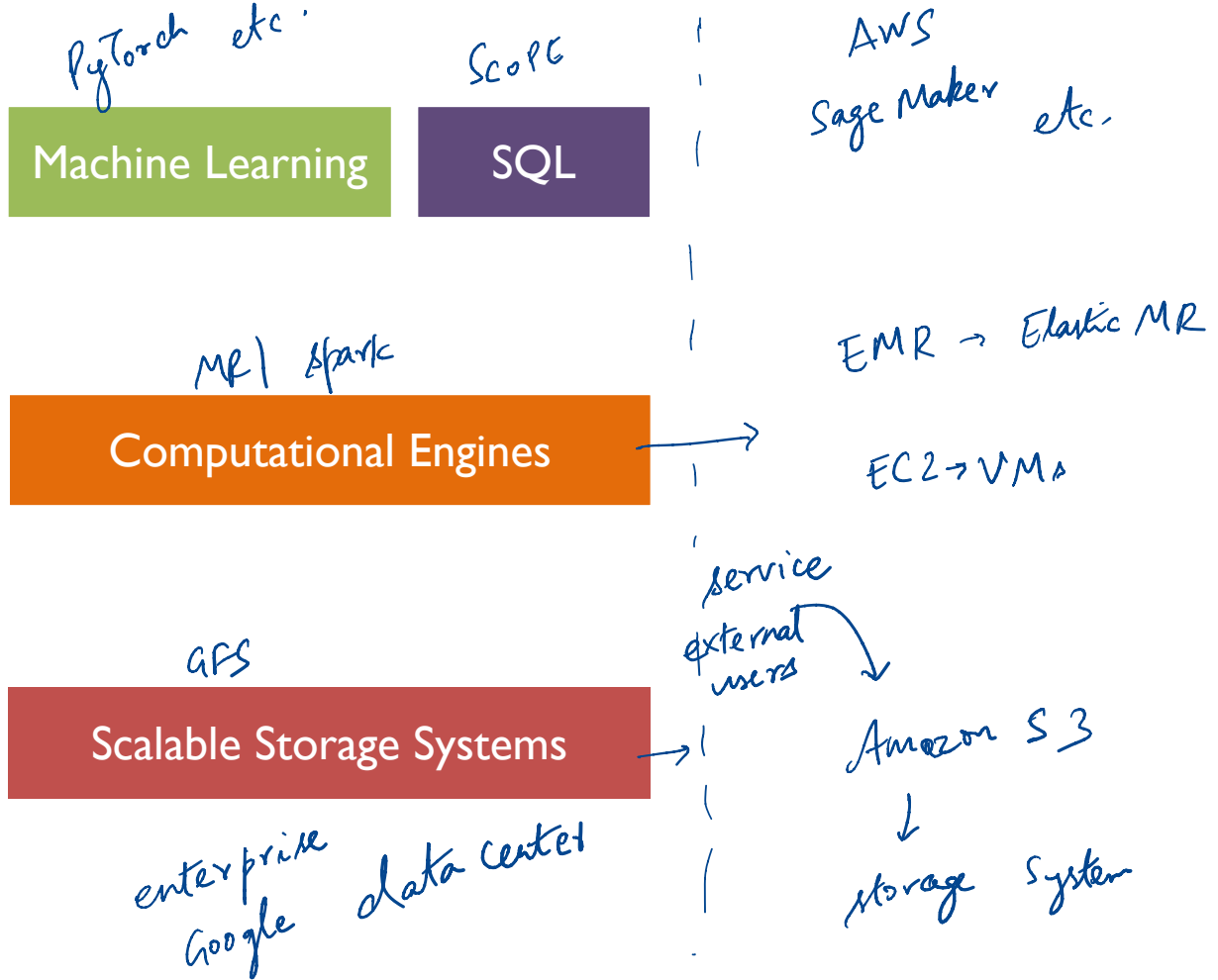
Snowflake

↳ elastic or
suited for the
cloud

CLOUD COMPUTING STACK

↓
Software as
a Service

↖ no management
overhead



SNOWFLAKE: GOALS

Software-as-a-Service → No need to download anything
Browser based & on-off as you need it

↖ Elastic → Use as many as required
Not any more / change when necessary
early

Highly Available → Fault tolerance

Semi-Structured Data → Data that might not fit a
strict schema vs. Relational
Databases
strict
schema

JSON,
XML
...

SNOWFLAKE DESIGN

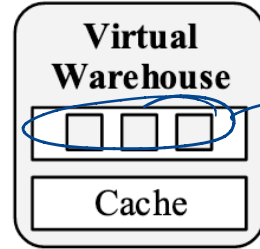
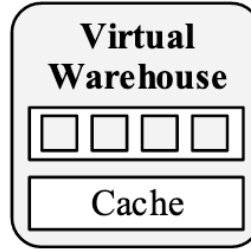
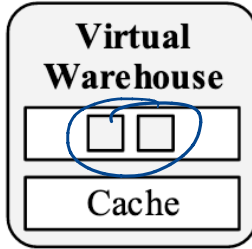
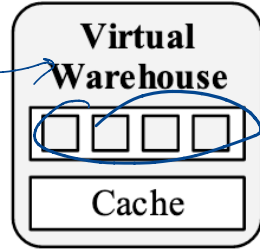
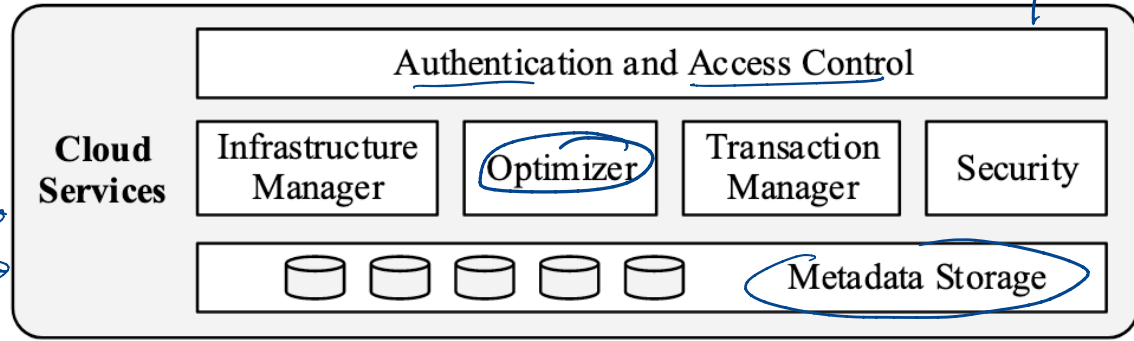
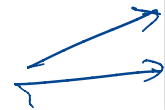
SaaS

Elastic

HA

Semi Structured

"Brains"
↓
SaaS
↳ multiple tenants / users



Elasticity
where each VW can have diff num machines

replicate to get HA

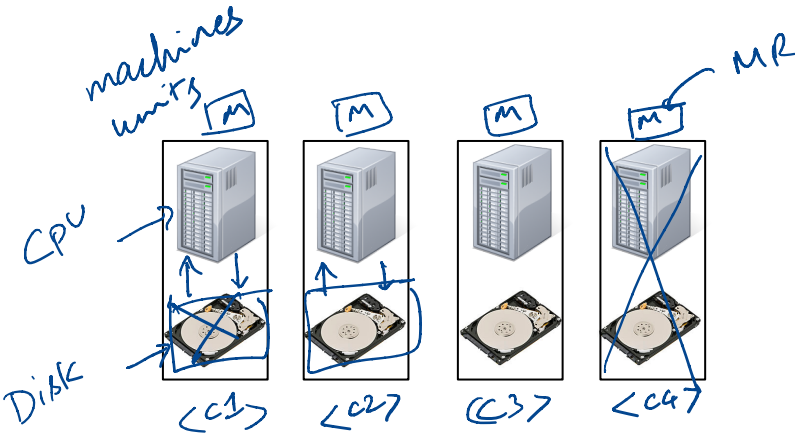


Share this data
↳ across jobs
↳ failures

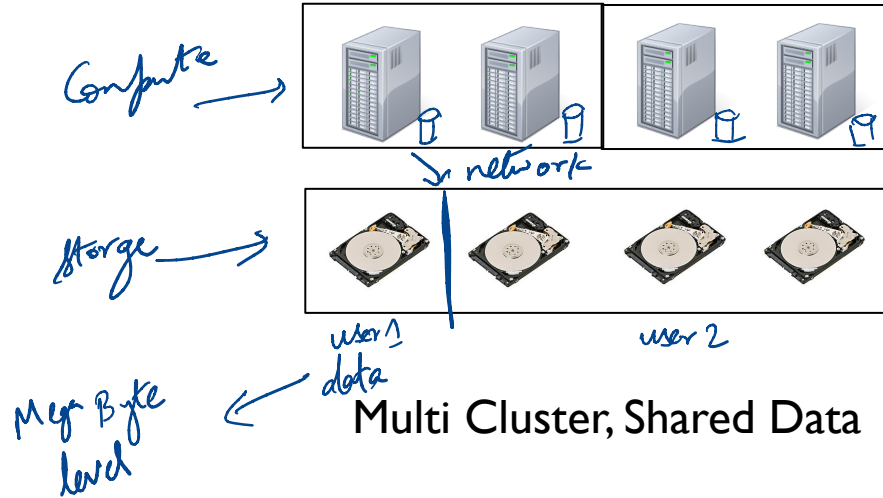
STORAGE VS COMPUTE

Coupling CPU, disk
↳ If I need 4 CPUs
⇒ 4 disks
↳ Utilization could be low

Separate out compute & storage
→ No locality!
→ Scale compute / storage independently



Shared Nothing
Each of these is a VM



Snowflake:
No database
indices!

STORAGE: HYBRID COLUMNAR

Table	Name	Age
rows	Alice	32
	Bob	22
	Eve	24
	Victor	27

horizontally partitioned
each partition
becomes a "file" in S3
↳ immutable

C1

Alice,32,Bob,22

C2

Eve,24,Victor,27

Row-oriented

Name
byte 0
byte 23
Age
byte 24
byte 32

name →
[Alice, Bob], [32, 22]

[Eve, Victor], 24, 27

Hybrid Columnar

A lot of queries
only touch a
few columns
⇒ range get
Avoid reading
entire
file!!

VIRTUAL WAREHOUSES

one user cannot slow down another user
↑

Elasticity, Isolation

↓
Every user can be diff in size

→ EC2 virtual machines
launched for a particular user
↳ Number of VMs → num, Type of VMs → cores → small → XXL

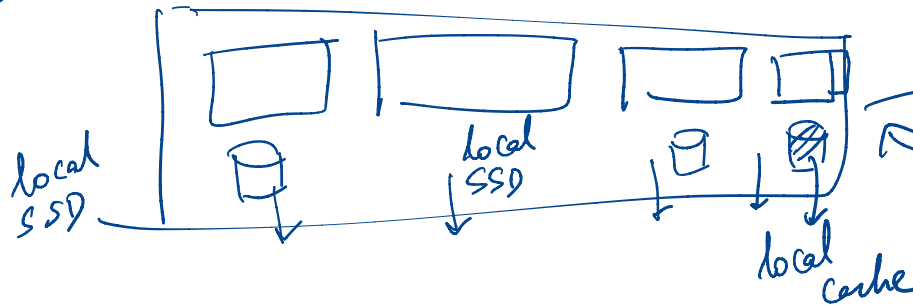
Only use them when running a query!

Local caching, Stragglers

simple LRU

≡ AFS / NFS but files are immutable

VW with 4 VMs

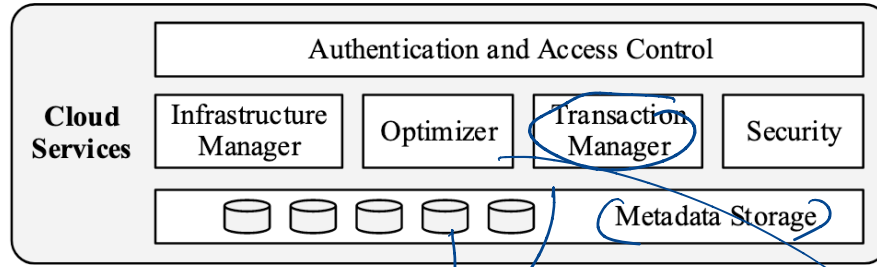


remote read S3

results

"File" stealing
→ if a worker is done early it steals work from slow worker

CLOUD SERVICES



query → Table V0
(c1, c2, c3, c4)

Table V1
(c1, c2', c3, c4)

↓
c2' → MVC C abort

Concurrency Control

↳ How to handle updates from many users

Snapshot Isolation = Isolation level

↓
all the reads come from a consistent version

Pruning

↳ Another way to improve locality

↳ It tries to skip files that don't have relevant tuples.

header Name / Age

	21	min 21
	32	max 32

← Age >= 35
Skip this file

FAULT TOLERANCE

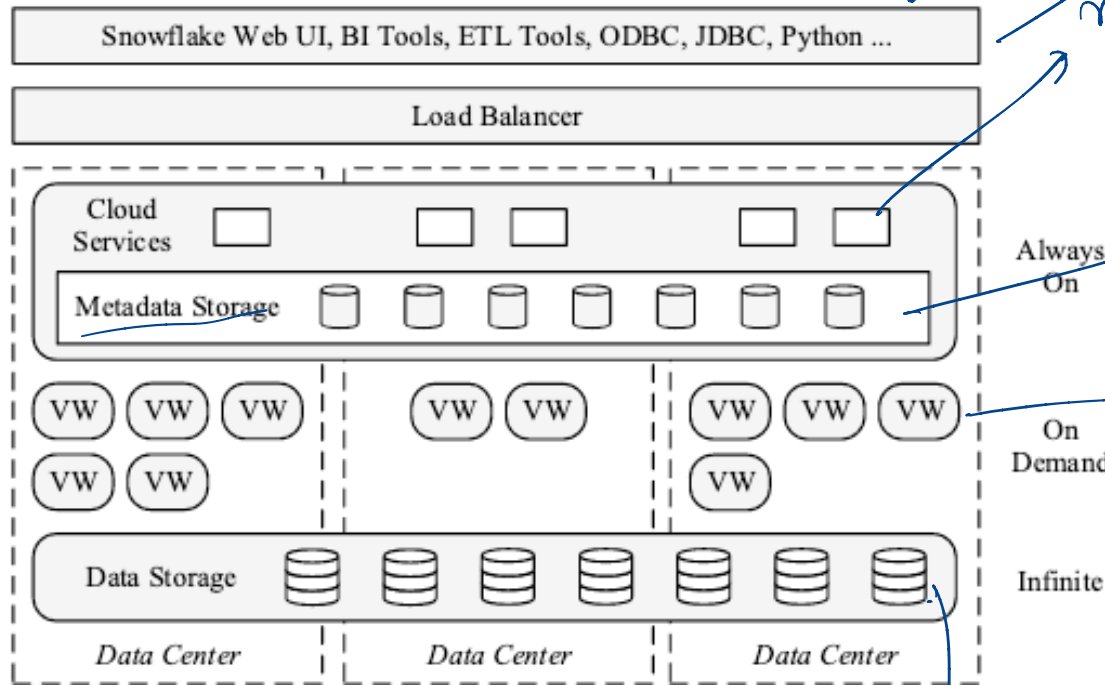
reminiscent of Ray GCS?
revert if services they fail

metadata

Replicate storage

Nothing. Ephemeral
If failure, retry the query.

Replicate data across data centers.



SEMI STRUCTURED DATA

```
{  
  first_name: "john",  
  last_name: "doe",  
  order_id: "1234",  
},  
{  
  first_name: "bucky",  
  last_name: "badger",  
  order_id: "52342",  
  order_date: "3/3/2020",  
}
```

no type
information

integer?

Some
fields

Extraction operation

↳ to access a field
inside JSON efficiently

Flattening

↳ arrays of JSON objects
can create rows out of
them.

Infer types, Pruning

Within a file

min
max

id	key	JSON blob
1234	1	...
52342		
...		

Serialized
JSON

within a file
try to infer
types

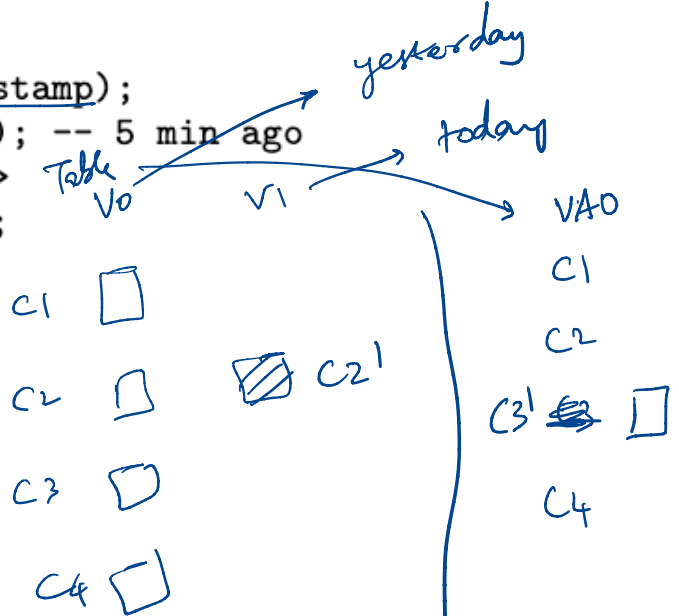
90 day retention policy

TIME TRAVEL?

→ over table versions

you can query multiple versions of a table

```
SELECT * FROM my_table AT(TIMESTAMP =>
'Mon, 01 May 2015 16:20:00 -0700'::timestamp);
SELECT * FROM my_table AT(OFFSET => -60*5); -- 5 min ago
SELECT * FROM my_table BEFORE(STATEMENT =>
'8e5d0ca9-005e-44e6-b858-a8f5b37c5726');
```

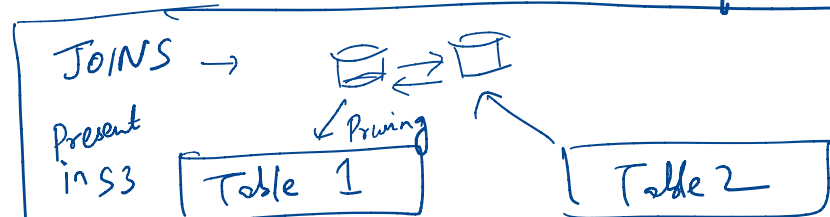


Multiple versions of table (MVCC)

Undo accidental deletes → new command **UNDROP**

Cheap to clone / snapshot a table

→ Copy on Write



SECURITY

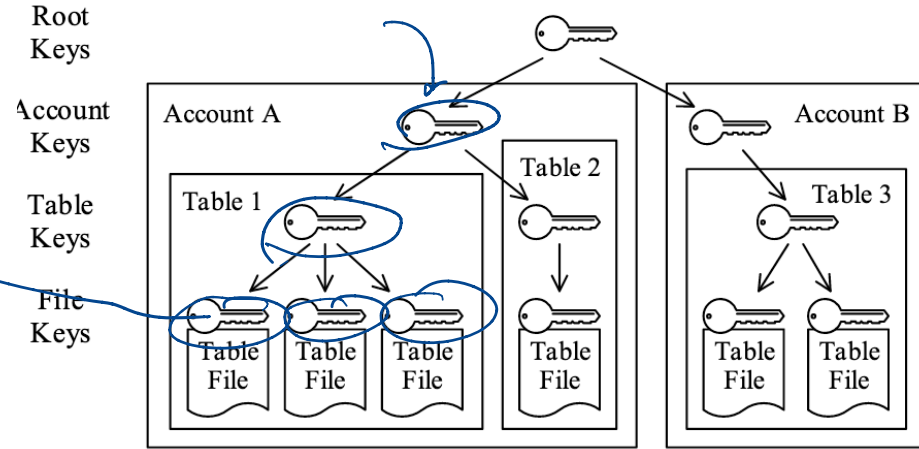
Hierarchical key management

↳ You encrypt a child key using parent key

→ limits what gets accessed when a key is compromised

Key rotation, re-keying

used refresh keys being used



SUMMARY, TAKEAWAYS

Snowflake

- Cloud computing → Elastic data warehouse
- Key idea: Separation of compute and storage!
- Hybrid columnar storage format → *rangers required*
- Elastic compute with virtual warehouses
- Pruning, semi-structured optimizations, fault tolerant
↓
limit number of files read

AEFIS FEEDBACK

DISCUSSION

<https://forms.gle/ZFosdUnizXYABAE86>

We see how Snowflake leads to the design of an elastic data warehouse. If we were to similarly design an Elastic PyTorch for training how would the design look? What are some design trade-offs compared to existing PyTorch?

version is created by a query

NEXT STEPS

Performance

Best performance within cost

Best cost within perf (5min)

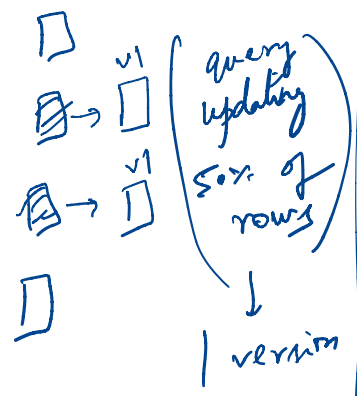
Next class: Midterm!

AEFIS feedback

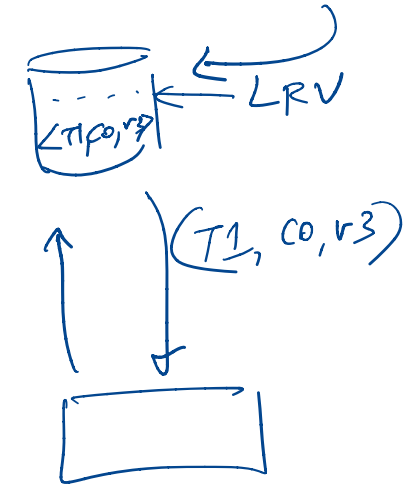
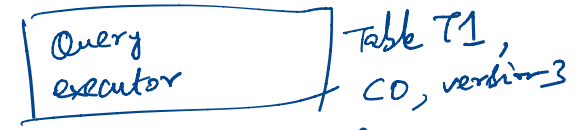
Project proposal peer feedback assignments

Indices
→ what is perf effect of not having them

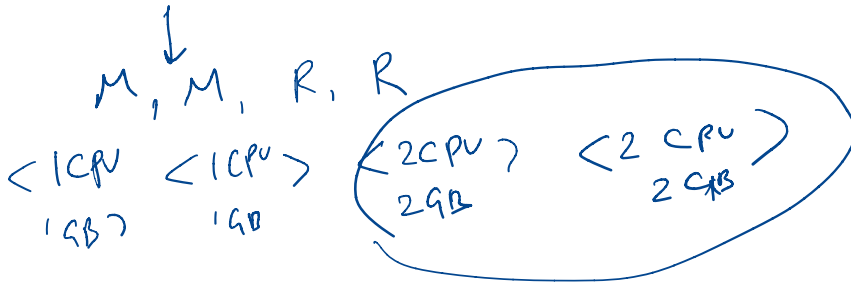
No need to submit discussion form!



Query → version of the table



DRF \rightarrow task dependencies



$\langle 6 \text{ CPU} \rangle$
 $\langle 6 \text{ GB} \rangle$ \leftarrow

Doesn't have a time dimension

Does this work?

DRF

→ Upstream tasks

Down stream tasks

Instantaneous
fair scheme

⇒ the downstream
tasks don't
inherit shares

RDDs

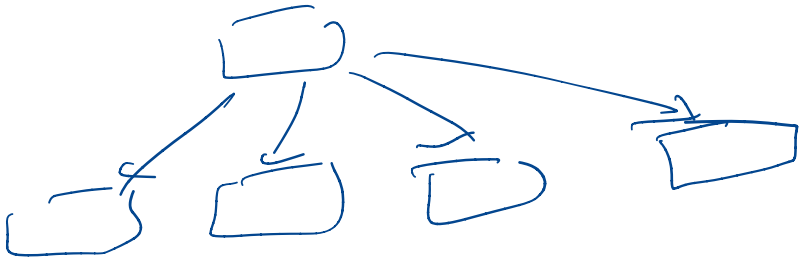
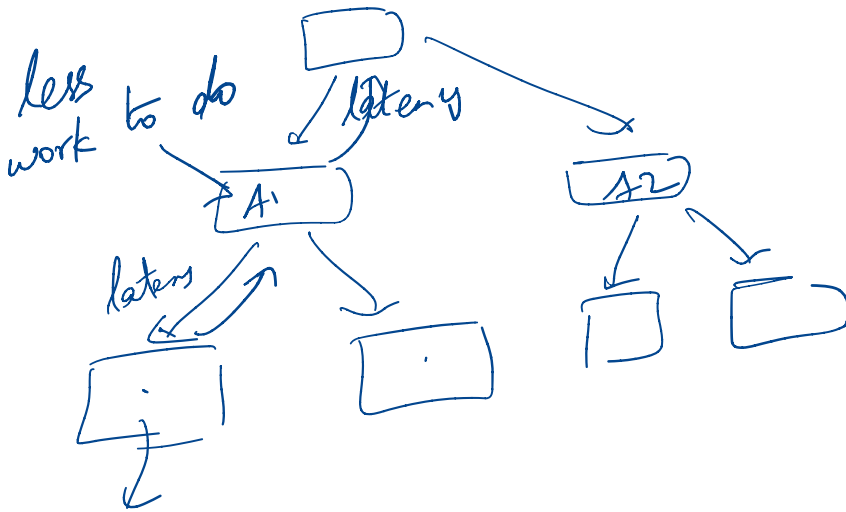
↳ Immutable vs. materialized

(b) Improving ~~FT~~ FT

↳ ^{just} lineage is default

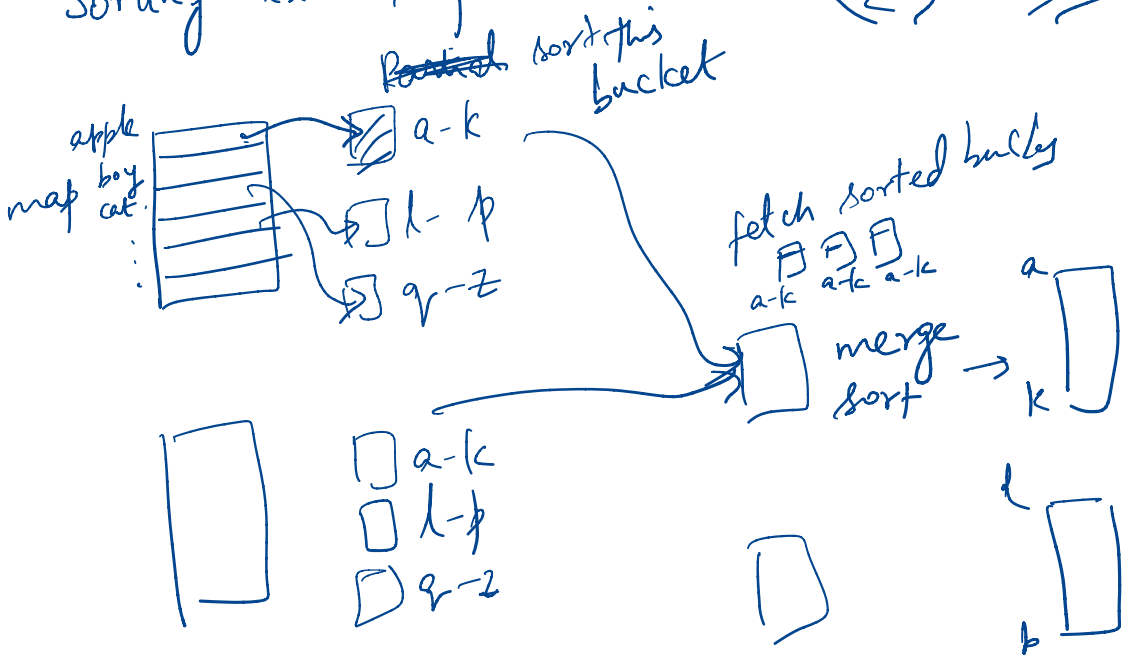
↳ ^{see} 5.4 in paper
checkpoints can shorten the
lineage

mid level aggregator



Sorting in MapReduce

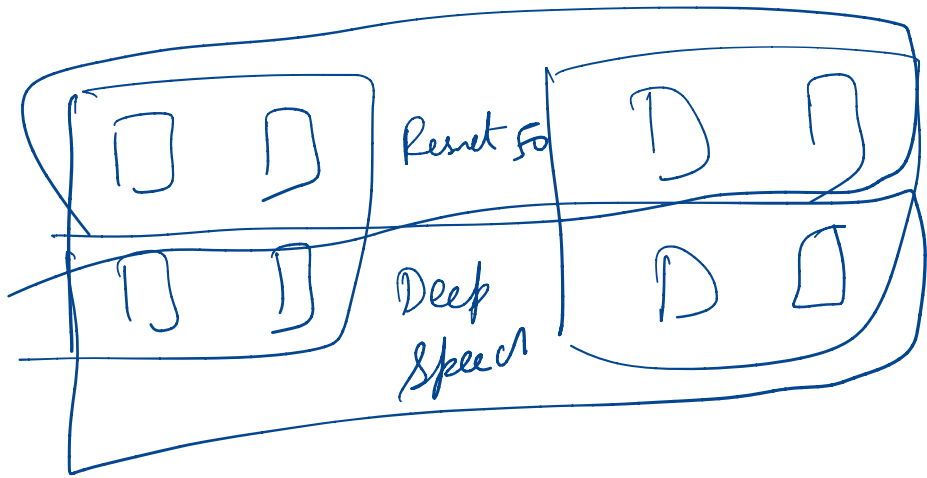
Stream



sample the data
random 1%
of data

histogram of words
fetch ~~it~~ to 1 machine
compute buckets

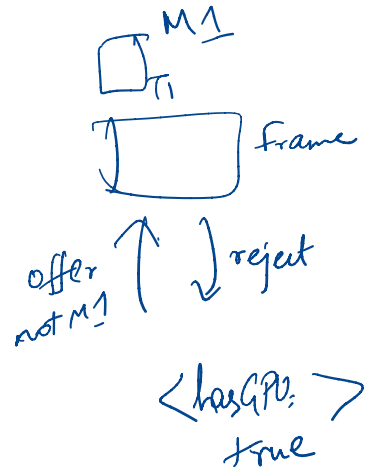
Gandiva, DRF



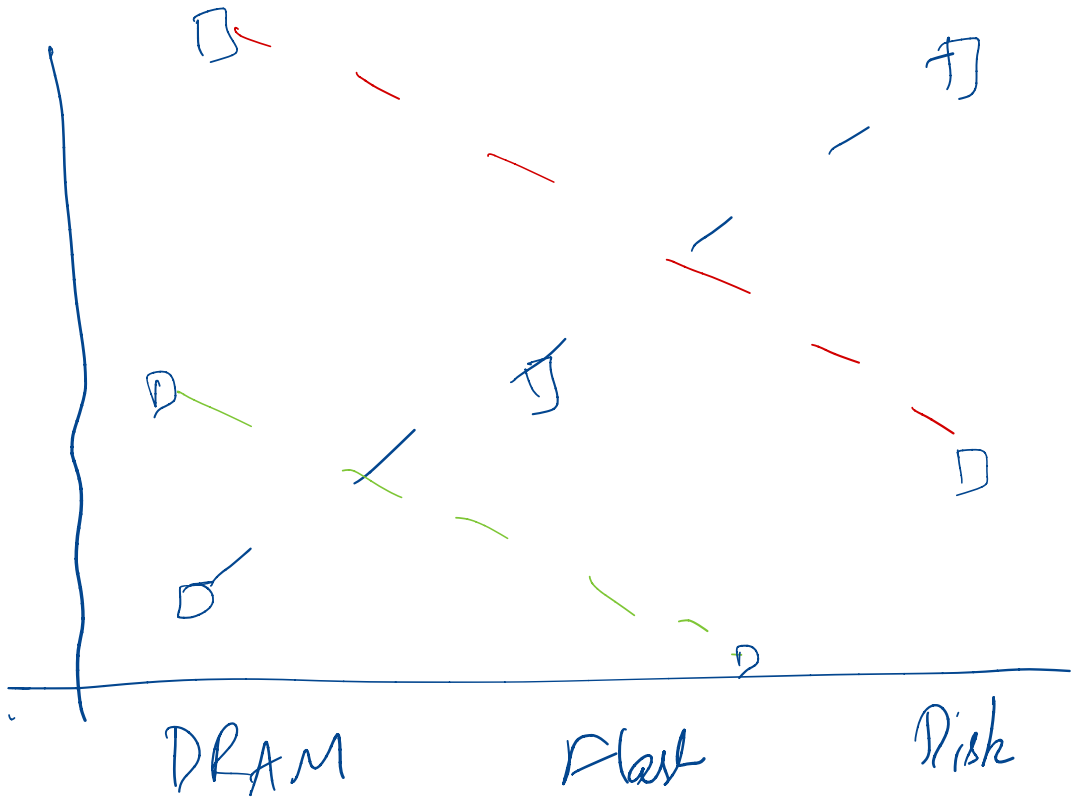
Sharing incentive if
shared allocation is
as good as having small
exclusive cluster

DRF \rightarrow task preferences

\downarrow
co-locate some tasks?



Workload



Blue: Capacity, latency
Red: Bandwidth, \rightarrow Power usage?
Green: Price / GB

MR failures

Assumption:
Process failure

- (i) Map complete failure
- (ii) In-progress map
- (iii) In-progress reduce

(a) Map output is already on disk, nothing to be done

(b) restart map task

(b) restart reduce task
all map outputs still available (only process failure)

