

1. OPENING

DATA LAKES

also

DATA WAREHOUSE

WITH

CLOUD

# # Information

## AA course series

1 diff data pw

2 data pipeline  
in batch

EC  
ECI  
EIL

code

3 streaming (pushes, delete (low)  
pg + blocks

4 short analysis  
(workflow, ML)

## Att course introduction

Spoke

- read of data set
- data and delete

## ## introduction to data eng

### ## modern info

a data eng build data pipeline

- challenge in this class
- big data solutions
- data lake / data diff
- perform with other tools
- data ethics policy and governance
- productivity and monitoring
- case study examples

### ## roles of data eng

data eng builds data pipeline

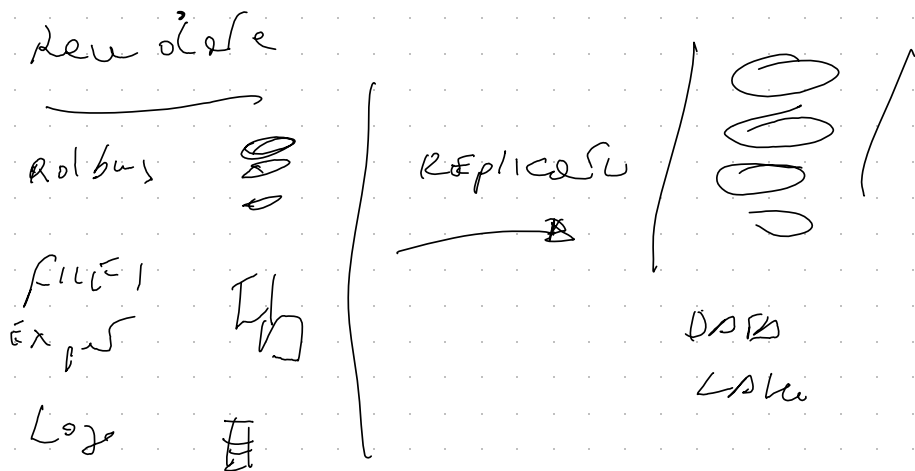
to have data in the place

to have data driven

business decisions

also a unstable condition

RAW state is not very useful



STORE RAW state also BUCKETS  
cloud split

think about

- state types to handle
- system fragments
- types of time - period  
access control

To make RAW state accessible  
for analytics

officially on this side you need  
ETL Extract Transform Load

using  
deprec or deflow  
streaming load  
using  
pub/sub or deflow

} Bp

the side of collection

- 1 - occurs to side
- 2 - side capacity and quality
- 3 - enough resources for Transf
- 4 - enough resource for purging side

① Ex data scattered across  
multiple systems / tools / schema  
commonly data is in files  
with each org having it's system

② Build ETL for ETL and Transform  
data

⇒ in EDW

data in EDW is

- clustered
- joinable
- predictable

③ on premise = FOL  
difficult — manage & provision  
scattered across  
files

Better JCP :)

# info to B2

$\left[ \begin{array}{c} B_1 \\ B_2 \end{array} \right], \sim$  physical  
merged  
CDW

per a map = jobs are organized  
into blocks

also like = def external schema  
as federated query

Tables and  
view = null sp

Creates = FSN for permission

12 Bp you don't need to provide  
!!

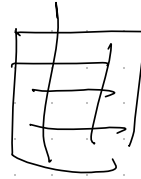
on-demand stored and output

query jobs / CPU + RAM  
25 days or

4th class lecture end - Du

1, 2, 3  
4, 5, 6  
7, 8, 9

REPLACES



ETL  
PIPELINE



Ren state

Defe  
late

Defe  
work hard

questions when building on

- Bosch early / early sink for pipelines
- 2 cells to meet needs
  - for some defe
  - for query performance
- few more is
  - organized
  - cell / per
  - e cells controlled

be 12 a modern DWH



you can use BP only as query engine

- cloud spl
- postgres
  - my sql

federated query

Cloud Storage

- CSB - FF
- join in
- persist

federated query

Bq

EQ



## ## Federated db vs dw

cloud spl is fully managed RDBMS

- automatic encryption
- 64 TB capacity
- 60 K R/W of 150ms
- auto scale and auto hdp

cloud spl is optimized for writes

Cloud

SQL

==

- scales to TB
- for backend ops
- RFS based stage

By

||

- scales to PB
- For analytics
- For ETL
- column-based storage

AFI partner with other teams

· NL Eng = 4m of BQ ML

Done Analysis

BQ <sup>or</sup> analysis

Done Eng

u

B ETL

(no autos)

Cloud monitoring to track queries  
and performance

## ## data access and persistence

we need a data provider model

- who should call it
- how PII are handled

also ↗

- data catalog → data discovery
- dlp → PII masking

## the prediction-ready pipeline

- ensures data pipeline health  
data cleanliness
- keep up to date
- handle errors

Apache AIRFLOW

is

data orchestrator

## How using BP for the analysis

- carry out queries
- combine all the analysis

### • Input

list of NOAA public water

### • Output

Open BP in cloud console

### - Task 1

Water project  $\rightarrow$  city bike - trips

### 1. Dataset

$min(start\_station\_id)$  is start-station

$max(end\_station\_id)$  is end-station

Approx. percentage (percentage, %)

$COFFSET(5)$  is typical-station  
from count (percentage) is non-typical

2. big-query-public-ny-city-bike-trips

WHERE

$start\_station\_id \neq end\_station\_id$

group by  
start\\_station\\_id, end\\_station\\_id

- task 2

o 1500

ux. del5,

ux. value / 10.0 es pcp

from

1. bicycle - public-del. gch-gz's  
es ux

where

1st 2 '05-121' and 2 flag in work

order by

ux. del5

- task 3

with bicycle - wheel ASL

1500

count (del5) as non-trip,

Extract (del5 from del5) as  
trip-del5

from 1. bicycle - public-del. trip

group by trip-del5

)

rainy - day, ex (

fields  
also,

(max(map) > 5) as rainy

from

(# sp/rock 2)

group by also

,

1/1/1/1

round (exp / blk, num - rps) as  
num - rps,

ex. rainy

from bicycle - rentals ex blk

join rainy - day, ex ex

or ex. day = blk. rps - day

order by ex. rainy

# BUILDING a DATA LAYER

## 1st introduction

- 1 - what is a DL
- 2 - data storage - places on PC
- 3 - build a layer of buckets
- 4 - faster access to blobs
- 5 - store all data files
- 6 - cloud API

## 2nd introduction to DL

by a scalable and secure and durable  
platform for compute - storage - analytics  
- everything

any type of volume of info

- standard - with - some
- back - many
- 171 - each - N/A

# EX Construction life

1 Role member  
assigned to  
Construction  
life  $\longleftrightarrow$  per defo  
info  
defect

2 member need  
to be cut  
out otherwise  $\longleftrightarrow$  pig, lotline  
info site sticks  
EDW

3 Building Floor  
and  
Roof  $\longleftrightarrow$  Tables, MC  
models,  
RF, PONT,

4 Supervisor  
all over  
and  
concrete life  $\longleftrightarrow$  work for  
architectural



#1 also stores the  $E_L$  options to keep about stages we know

- cost stage
- cost sp
- cost power
- Filter
- + - cost by factor

Keep in mind:

- where else is now
- how big is
- where else has to go a desk
- How much to transfer in  $\Delta t$

Consider

$E_L =$  data input in a format compatible to load in desk

input also into bp

ELT = we need to  
check in  
transformation  
✓ we can

Load from source is about  
now then apply T logic

use bp to write up to transformation  
input side and write it into  
now table

ETL = we apply T before  
to load

T greatly reduces the size

binary proprietary input format

T it before to load into target

usually, dataflow pipeline  
before  
writing into bp

~~It build also takes way CS~~

CS storage characteristics

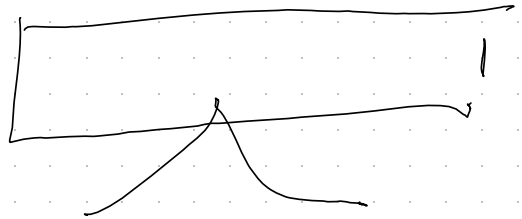
- persistence
- durability
- strong consistency
- availability (global on super zone)
- high throughput

It is a object store

with low filesystem structure

a Bucket contains objects

Bucket



objects



bucket is globally unique  
associated with region

low express cheap, cheap  
low region

multi region bucket

↳ replicate across regions

single region bucket

↳ replicate across zones

metadata is stored along the obj

obj

obj

obj metadata

compression  
encryption  
lifecycle mg

Storage class

- frequent
- non-frequent storage
- Cold line storage
- Archive storage

1 30 days  
once a month  
4 90 days  
once a quarter  
4 365 days

# #11 is AWS cloud steps

control access to S3

IAM

(1)

is a gap  
for all  
product

ACL

(2)

apply of bucket  
to S3

finer

bucket  
IAM roles =

- Reader
- Writer
- Owner
- (if ACL policy)

delete/control bucket is a S3  
level role

both read/write bucket roles  
access to policies w/ ACL

Encryption is now in

GREEN

DEK + key encryption key

by Google rotated on schedule

CMEN

by customer myself

C>EN

by customer implied

+ client-side encryption  
11 extra-options

You can set lock policy to 96  
bucket cannot be changed

## first element of state

Transactional system are 30%, 40%, 20%, 25%

Analytical systems are populated from transactional system

if you require low latency  
in order of milliseconds

↓  
cloud byt

## cloud is relational database

cloud sql

- fully managed

(2021) db services

- sql database

- my sql

- postgres sql

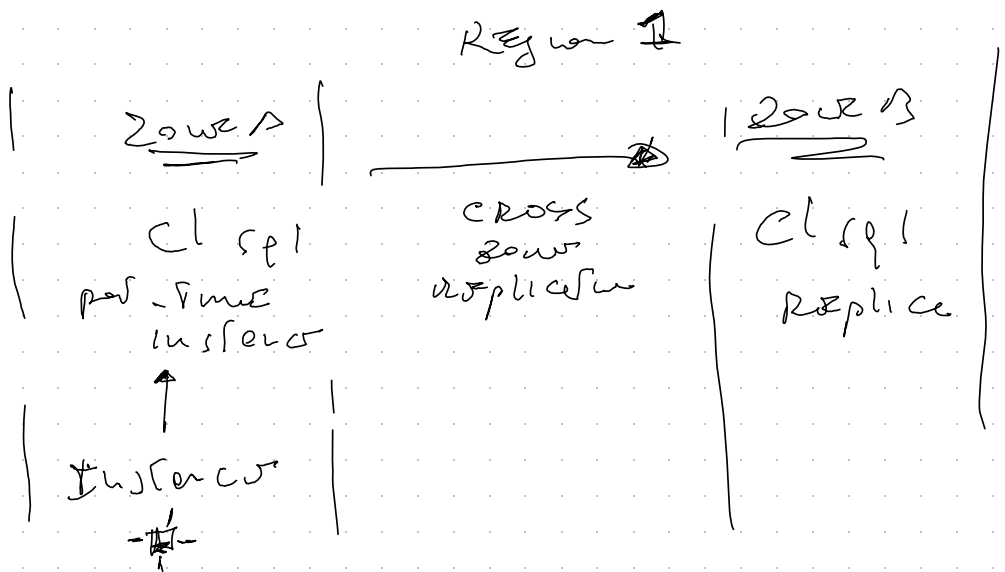
scale vertically (R+U)

- 64 process
- 100 GB RAM

scale horizontally (R)

- Replicas

> if you need HORIZ (R+U)  
consider Cloud Provider





Modul 15/16/17, alle auf 20C

2501 Tue

Best and

pub / sub Async msg

Asymptotic

1  
be-flow

Charles  
Henry

Clear By ✓

Spring | | 2015  
pipe/w pipe/w

has a  
pipeline

Bf

Delema

# Muscular Exploration

- Leoben
- Talsu
- Spreitzer

#11 Les and For. store w/ global  
if,

- create cloud spl instance
- create new obj
- input
- check data integrity

• For 1

0 Export project\_ID = \$ (global info  
-- from = 'value (config. project)')

1 Export BUCKET = \$ { project\_ID } - .ml/

• For 2

0 global spl instances created  
Taxi -- Tier = alb - ml - s3 - 1  
-- execution-policy = ALLOW

1 global spl users if - prod test

2 # ip addresses

Export ADDRESS = \$ (ipof - prod - )

3 global ——— patch taxi = execution  
networks,

• look >

7 gcloud upgrade --project \$PROJECT\_ID --version \$VERSION

1) mysql --host=\$MYSQL\_HOST --user=\$MYSQL\_USER

2) Load DATA

# BUILDING A DATA MARCH

## #1 Architecture

- 1 modern architecture
- 2 info to bp
- 3 lead olek in bp
- 4 target scheme
- 5 scheme only
- 6 horizontal repeated field

## #2 The modern olek network

- <sup>scout from</sup> sig byte  $\rightarrow$  Pb
- universal (w-o-p, l)
- highly eco-system
- etc eco-system
- up to network level
- 100% available with no olek now
- security and collaboration

## unlike to be

is similar to cloud layer  
to no rooms

GIS and ML

clustering data inside one epd

ANALYZE DATA

no vacuum power, like resistor  
no IDA to be deleted

by itself now alone executed

we have stage engine  
and

compute engine

no need to provision of resource

static unit of computation  $\begin{cases} \text{cpu} \\ \text{mem} \\ \text{network} \\ \text{storage} \end{cases}$   
any 2 or 3 for user

#The gif started with bp

1) תוספת פחם ופחמן דו-חמצני

"פשוט . מפורסם . ידוע"

de 10 155 = electrolyzed water  
proj's never copy to water occur,  
ely biller  $\rightarrow$  project

you can e query in yellow prj  
even use sfqa defc in another  
project

for each guess  $c$  → I AM able to submit a job

AC is the offset or position of  
V<sub>max</sub> or column level

bp. steps can be — replaced  
— multi-replaced

each table has a refresh  
choice

query solutions can be used  
or  
payd. work/poc accounts

Admin & system events are  
all logged

table expiration ~ system event

For predefined roles

- hp admin
  - hp role expires
  - hp role summary
  - hp role statement
  - hp job user
  - hp user table user
- current role/permissions  
with views

You can select — authorized  
views

→ role/permissions  
permissions

CREATE ROW ACCESS POLICY

spec-filter on [table]

GAINS to ("group: sales - spec  
example - com")

FILTER USING

is (Region = "APAC")

	person	colours	country	Region
dc				APAC
no				US

Authorized view →

user can access view

but not

underlying tables / data

diff mechanism view ————  
better  
reference  
Δ store  
local as  
autonomously



4# lead able into bp

the hot 5C / 5C / 5C / 5C

you can watch lead able into bp

CLV - 100% - one, perfect, one  
gap compressed, first four years

Lead job - created for of table  
if not exist

1 chance is given of lead job

be has child on two of jobs

and still you can lead every day

Be a 100% mostly by

overflow - overflow

be has connections

Be mention a 100% history

CREATOR DESCRIPTION

○ CREATOR OR REPLY (color  $\Gamma_1$ )  
A)

UPPER & FLOW ( — )

FOR SYSTEM TIME AS OF

TIME(S) (SUB/CURRENT- $\Gamma_1$ ,  
INTERNAL 244)

By transferable service also  
automatic also transfer

Be Transfer some provide,

- connections (lost)
- synchronization problems
- scheduling

o global storage of x.c.v  
js: // my bucket

! be load —

be inputs delivered are

INSERT / UPDATE /  
DELETE / MOVE

! be is not ok

help to create complex expression  
in another helper

o create temp function multiply(x,y)  
returns result

LONGER JS AS more return x \* y; u

# body here ~ bp : 43

- body obj. < n bp
- than - ——— my in / context
- my obj

look 3

bp head -- source-func = CPU \

-- end of file \

-- next place \

obj. file. table A

go: // cloud-ky / obj / 129. CPU

# explor schemes

information - scheme ——— TABLES

——— JOURNAL

—— TABLES ——

#4 scheme strategy -

by store k

customer	store id	date	items
Doug	123	2021/01	Ref, Power
Alex	124	—	1kg, orange potato, 10kg Cipolla

Normalized store  
products

A1	123	2021/01
A2	124	—

customer

Doug	A1
Alex	A2

store - items

123	Ref	—
123	Power	—
123	orange	1kg
124	potato	—
124	Cipolla 10kg	—

rows col by row

by column

A

— have  
— price

diff way to represent data / Fast  
ACG

whether you demonstrated skills  
before getting into by

Be specific columns with  
preferred and  
preferred date

<u>order ID</u>	<u>order date</u>	<u>ord. product</u>	<u>ord. quantity</u>
123	20/10	person	1
		coat	1
		people	10
124	✓	version	
		orange	2

In This way he can experience  
the date in some way / collected  
still

May a file/demonstrated structure

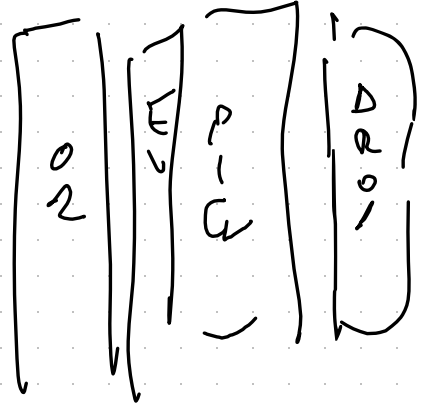
# ## noted and repeated fields

order

events

Pickup

Drop  
off



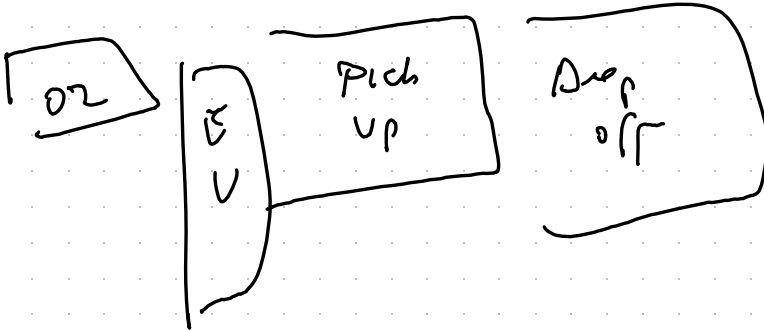
many jobs

2 by 10s

early joins

also 11  
repeated

✓ reference



rec used and // diff granularity  
repeated

# ## working with floor lab

- body maintenance job
- closely and purely of array  
↓  
fruits
- punning up harvest  
↓  
2 specified fields

Task . 1

create array → fruit - 1025

Task . 2

row	Fruit	Price
1	apple	1000
2	banana	1000
3	cherry	1000
4	orange	1000

specified fields!

row	fruit (array)	Price
1	(apple, banana, orange)	1000
2	(cherry)	1000



array only 1/4 E

↓ then

new vector fruit-ions

↓ now

For 3

ARRAY - AGG)

↓ SELECT

full visitor EA,

defo,

ARRAY - AGG (product-line) as p4

ARRAY - AGG (prop-5112) es p5

FROM - e-commerce. de-headers

group by full visitor EA, defo

RECAP

ARRAY - CEMENT (cement)

ARRAY - AGG (DISTINCT (field))

ARRAY - AGG (C field) ORDER BY C - 1)

ARRAY - AGG (C field) LIMIT N)

Control access, a field present on  
a volume with type character,

! before we can perform operations  
you must first know source -> rows

hits - page - page title

fewer distinct  
visited,

4-univers page - page title

from 'pe-headers'

on most (hits) as 4-univers

links

Top < structure

- 4 + fields
- some on diff types
- other structure

links -

total \* , done \* -

Task 6.

jQuery ( "body" ) .html ( "

23.5 <div> 23.2 <div> 20.1 <div> 20.1 </div>

" )

jQuery ( "div" ) .html ( "

23.4 <div> 23.2 <div> 20.1 <div> 20.1 </div>

" )

source as text:

{ type: "text", value: "23.5 <div> 23.2 <div> 20.1 <div> 20.1 </div>" }

fields:

{ type: "text", value: "23.5 <div> 23.2 <div> 20.1 <div> 20.1 </div>" }

{ type: "text", value: "23.5 <div> 23.2 <div> 20.1 <div> 20.1 </div>" }

Solve solution - be careful

Field Name?	Type	Notes
RACE	STRING	NULLABLE
♦ PARTICIPANTS	RECORD	REPEATED
NAME	STRING	NULL
Split	FLOAT	REPEITION

Section - From Recy. use results

row	recy	participants name	performance split
1	900m	Rudisha	26.3
			25.2
			24.1
		Dehaish	24.3
			' '
		Murphy	23.1
			24.2
			0.0

SELECT recs, participant names

FROM recip. recs - results

CROSS JOIN

~~participants~~ \* study - active  
table inside the  
table

recs - results

# // full name

participants

or platform

SELECT recs, participants names

FROM recip. recs - results AS r

,

r. participants

Table 1.

is group

SELECT count (p. names) as pc

FROM recip. recs - results AS r

,

UNION (r. participants) AS p

Task 3.

10155

p-norm

$\text{AVG}(\text{split-func})$  is  $\text{avg-f}$

From recurf. case, results as 2

$\text{UNNORM}(\text{2. probecount})$  is  $p$

$\text{UNNORM}(\text{p-split})$  is  $\text{split-func}$

group by p-norm

order by avg-f desc

## optimizing with partition and clustering

reduce cost with partition

C1 C2 C3      evaluate      C4

↓  
data is split in separate threads  
parallel then will be ready

3 way partition

- 14 partition size -

by party -- observation - schedule d.f

-- size - partitioning - type = day

- on cd of type also time, data,  
time stamp

by mk -- table -- schema

o: string

b: float

-- time partitioning. first

- ungroup type  
columns

-- range partitioning =  $\left( \begin{array}{l} \text{customer id} \\ 0, \\ 100, \\ 10 \end{array} \right)$

know our sql performance

SELECT

Act,

FROM

myowerf. table

WHERE

- PARTITION TIME >



TIME RANGE 100

( TIME RANGE (

date

on

LEFT most

1 2022-07-01

INTEGRAL 5 DAY

the partition field

by group -

-- requires partition filter



clustering on field helps too

CREATE TABLE mytbl. mytbl (

col1 VARCHAR,

col2 VARCHAR,

col3 VARCHAR,

col4 VARCHAR FIRST LAST

)

PARTITION BY DATE (col4)

CLUSTER BY col1

options

(partition-requests=1,)

AS

tbltbl & FROM mytbl mytbl

He does also RE-clustering

for you

use clustering

your table is partitioned  
on col4 / last

use agg. / filter on  
column / fields

