

Udacity : Data Engineering

1. Welcome

03/18/2021

Project 1: Data Modeling

- PostgreSQL, NoSQL - Apache Cassandra

Project 2: Cloud Data Warehousing

- ELT pipeline

- S3, Redshift

Project 3: Data Lakes w Apache Spark

- ELT from S3

- process w Spark

,

Project 4: Data Pipelines w Apache Airflow

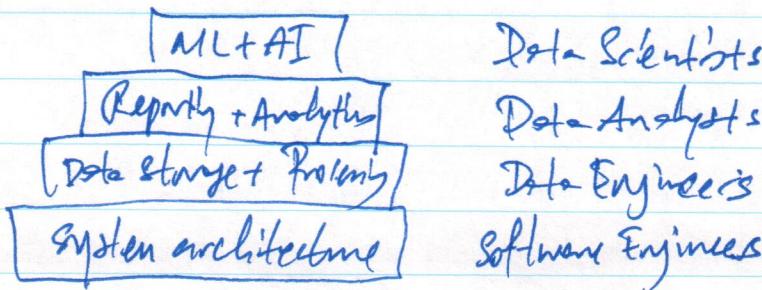
- scheduling pipelines + triggers

Project 5: Data Engineering Capstone

2. Data Modeling

2.1 Introduction to Data Engineering

Data engineering comprises of all engineering + operational tasks designed to make data available for end-user for analytics, model building, app development, ... etc.



Data Scientists + ML Engineers

Data Analysts

Data Engineers

Software Engineers

Common DE duties

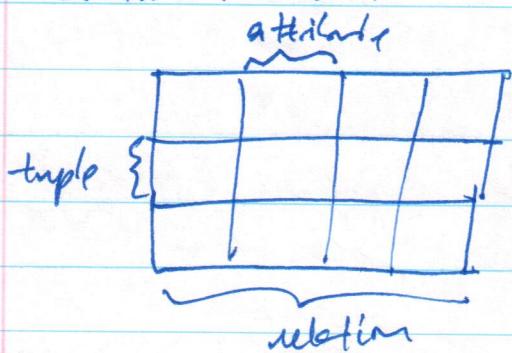
- ingest data from source
- build + maintain data warehouse
- create pipeline
- create analytical tools
- migrate data to cloud
- schedule + automate pipelines
- backfill data
- optimize queries
- design database

2.1 - Data Modeling

Data Modeling - organize data and how they relate to one another

- conceptual data modeling
 - relationship between tables
- logical data modeling - connect tables DDL
- physical data modeling (data definition language)
 - actual coding/implementation

Relational Model



SQL -

- Oracle
- Teradata
- MySQL
- PostgreSQL
- SQLite

Database / Schema - collection of tables

Table / Relation - groups of rows (tuples, attributes)

Advantages of Relational Databases

- ease of SQL
- JOINS
- aggregation + analytics
- smaller data volumes
- easier to change based on new requirements
- flexibility for queries
- modeling data, not modeling queries
- secondary indexes available
- ACID transactions - data integrity even in
 - guaranteed validity, even of commit, power failure
- Atomicity - all or nothing
- Consistency - have to abide by rules/constraints
- Isolation - processed independently + serially, order does not matter
- Durability - completed transaction are ^{committed} saved, even if power failure.

When not to use Relational Database

- large data - RDBMS can only scale vertically
- different data types
- high throughput - fast read
- flexible schema
- high availability - little downtime
- scale horizontally - add more servers to increase performance

PostgreSQL

- open source, SQL

```
connection = psycopg2.connect("host=127.0.0.1 dbname=studentdb  
user=kay")
```

```
cursor = connection.cursor()
```

```
connection.set_session(autocommit=True)
```

```
cursor.execute("create database music_library")
```

```
cursor.execute("Create Table If Not Exists music_library")
```

```
(album_name varchar, artist_name varchar,  
year int);")
```

```
cursor.execute("select count(*) from music_library")
```

```
print(cursor.fetchall())
```

```
cursor.execute("Insert Into music_library (album_name, artist_name,  
year)  
Values ('7s', '6s', '7s')")
```

```
["    ", "    ", "    ", 1970])
```

```
cursor.execute("Select * from music_library;")
```

NoSQL

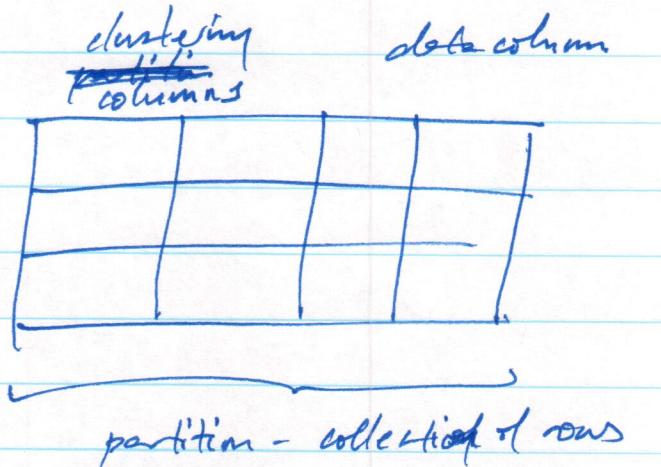
NoSQL = Not Only SQL

horizontal scaling

- Apache Cassandra (Partition Row Store)
- MongoDB (Document store)
- DynamoDB (key-value store)
- Apache HBase (Wide column store)
- Neo4J (Graph Database)

Apache Cassandra

- keyspace - collection of tables
- table - group of partitions
 - row
 - does not allow duplicates
 - scalability
 - high availability / fault-tolerance
 - CQL query language
 - good for - time series logging
 - IoT
 - time series data
 - any workload that writes to database
 - can't do ad-hoc queries, but can add clustering columns + create new tables



partition - collection of rows

When to use NoSQL

- large data
- horizontal scalability
- high throughput - fast reads
- flexible schema
- high availability
- store different data type
- user one distributed - low latency

When not to use NoSQL

- need ACID transactions
- need to do JOINs
- need to do aggregation + analytics
- changing business requirements
- queries not available, need flexibility
- small dataset

2.2. Relational Data Models

Database Management System - DMS

- all information represented by tables

Relational importance

- standardize of data model
- flexibility of adding + altering tables
- data integrity
- SQL
- simplicity
- intuitive organization

Online Analytical Processing (OLAP)

- optimized for reads, allow for complex, ad hoc queries
- heavy on aggregate

Online Transactional Processing (OLTP)

- allow for less complex queries in large volume
- optimized for read, insert, update, delete
- little aggregate

Structuring Database

- Normalize - reduce data redundancy
 - increase data integrity
- Denormalize - done in read heavy workloads to increase performance

Objectives of Normal Form

- free database from unwanted insertion, update + delete dependencies
- reduce need for refactoring database as new types of data are introduced

- make relational model more informative to users
- make database neutral to query statistics

First normal form (1NF)

- atomic values - each cell contains unique and single value
- able to add data without altering tables
- separate different relations into different tables
- keep relationships between tables together w/ foreign keys

Second normal form (2NF)

- reach ~~1NF~~ 1NF
- all columns in table must rely on Primary key

Third normal form (3NF)

- reach 2NF
- no transitive dependencies

When update data, just do in 1 place.

e.g.	Album ID	Album Name	Artist Name	Year	list of Songs
	1	Rubber Soul	Beatles	1965	[, ,]
	2	Let It Be	Beatles	1970	[, ,]

1NF:	Album_ID	Album_Name	Artist_Name	Year	Song_Name
	1	Rubber Soul	Beatles	1965	:
	1	"	"	"	:
	2	Let It Be	"	1970	:
	2	"	"	"	:

2NF: Album ID	Album Name	Artist Name	Year
1	Rubber Soul	The Beatles	1965
2	Let It Be	"	1970

Song ID	Album Name	Song Name
1	1	-
2	1	-
3	2	-
4	2	-

3NF: Song ID	Album ID	Song Name
1	1	-
2	1	-
3	2	-
4	2	-

Song library

Album ID	Album Name	Artist ID	Year
1	Rubber Soul	1	1965
2	Let it Be	1	1970

Album library

Artist ID	Artist Name
1	The Beatles

Artist library

Denormalization - joins are slow

- improve read performance of a database at expense of losing write performance by adding redundant copies of data
- common after normalization
- need more space due to copies
- write slower + on multiple tables

Fact Table - measurements, metrics or facts

Dimension Table - people, products, place + time

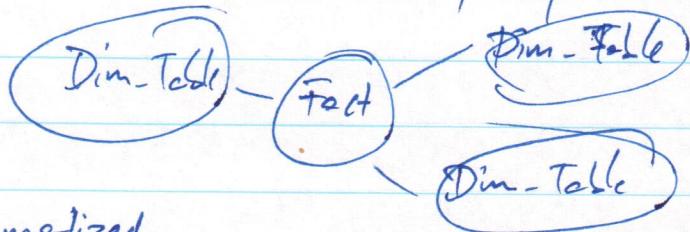
Star Schema

Snowflake Schema

} Entity Relationship Diagrams (ERD)

Star Schema

- > 1 fact table referencing any dimension tables



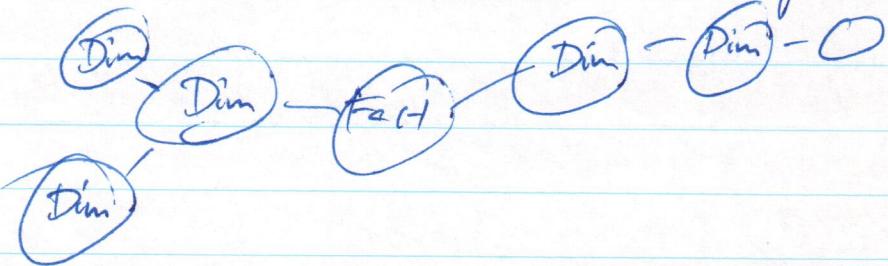
- denormalized
- simplified queries
- fast aggregation

Drawbacks

- issues from denormalization
- data integrity
- decreases query flexibility
- many-to-many relationships

Snowflake Schema

- centralized fact table connected to multiple dimension tables.



- does not allow for 1 to many relationships, while star does
- more normalized than star, but only in 1NF or 2NF

upsert

2.3 Project 1 : Data modeling in Postgres

- fact + dimension tables for star schema
- ETL pipeline to transfer data from files in 2 local directories into tables in Postgres
- Song dataset - JSON - metadata
 - num_songs
 - artist_id
 - artist_latitude, longitude, location, name
 - song_id, title, duration, year

- log dataset - JSON - year/month/events.json
 - artist
 - auth
 - first_name, gender, last_name
 - item_in_session, length, level, location, method, page, registration
 - session_id, song_id, status, ts, user_agent, user_id

- fact table - song plays - songplay_id, start_time, user_id, level, song_id, artist_id, session_id, location, user_agent

- Dimension tables - users - user_id, first_name, last_name, gender, level

- songs
- artists, - time

2.4 NoSQL Data Models

Not Only SQL

Apache Cassandra - open source NoSQL

- high availability
- linear scalability
- created Facebook
- masterless architecture
- (node = server, system, ...) each node connected to another node.
- users - Netflix, Apple, FB, Twitter

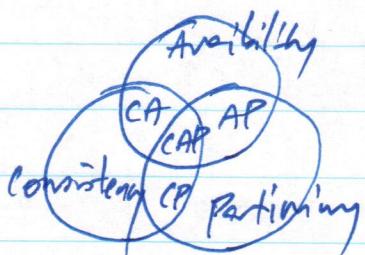
Distributed databases - for high availability, need copies of data.

Eventual consistency - if ^{no} new updates, all data in nodes are same.
but for few seconds of latency, maybe inconsistent.

CAP theorem - impossible for distributed ~~systems~~ data store to simultaneously provide 2 out of 3 guarantees

- consistency, availability, partition tolerance

every read /
give correct data always on
)
) work regardless of
) loss of network



C in CAP is not ACID
every read
has correct data.

Transaction must
be valid by constraints +
rules.

mostly C and P are guaranteed.

Denormalize in Apache Cassandra

- no JOINs.
- consider query first/before designing tables.
- ~~optimize~~ partition for fast reads.
- 1 table per query is good strategy, allows for redundancy.
- must be denormalized

CQL - similar to SQL, but JOINs, GROUP BY

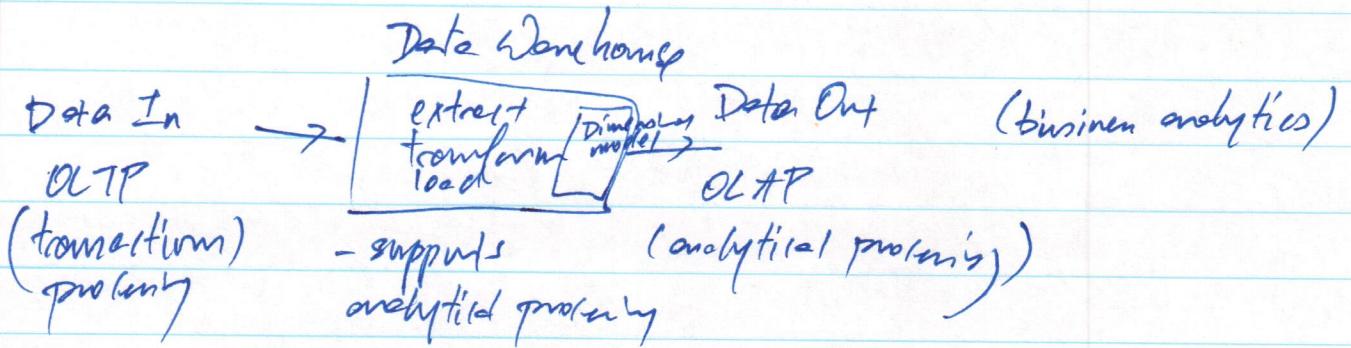
Primary key - uniquely identified, no duplicates clustering /
- 1st element is Partition Key or plus ~~Partiti~~ column
- Partition key is hashed into numeric
- data distributed by this partition key
- simple or composite
- ≥ 1 clustering columns

Primary key (a) - partition key \Rightarrow a
" " (a, b) - " " is a, clustering key \Rightarrow b
" " ((a, b)) - composite partition \Rightarrow (a, b)
" " (a, b, c) - partition key is a, composite clustering key (b, c)
" " ((a, b), c) - composite partition \Rightarrow (a, b), clustering key \Rightarrow c
" " ((a, b), c, d) - " " " (a, b), " " " \Rightarrow (c, d)

3. Cloud Data Warehouses

3.1 Introduction to Data Warehouses

- Operational process - "make it work"
- Analytical process - "what's going on"



Dimensional Modeling (Rep)

- star schema - good for OLAP, not OLTP
- Normalized Farm (3NF)
- Fact, Dimension tables - dates + time, locations, values
(facts are numeric + additive (usually))

% load-ext sql

% sql select count(*) from table; \Rightarrow single line

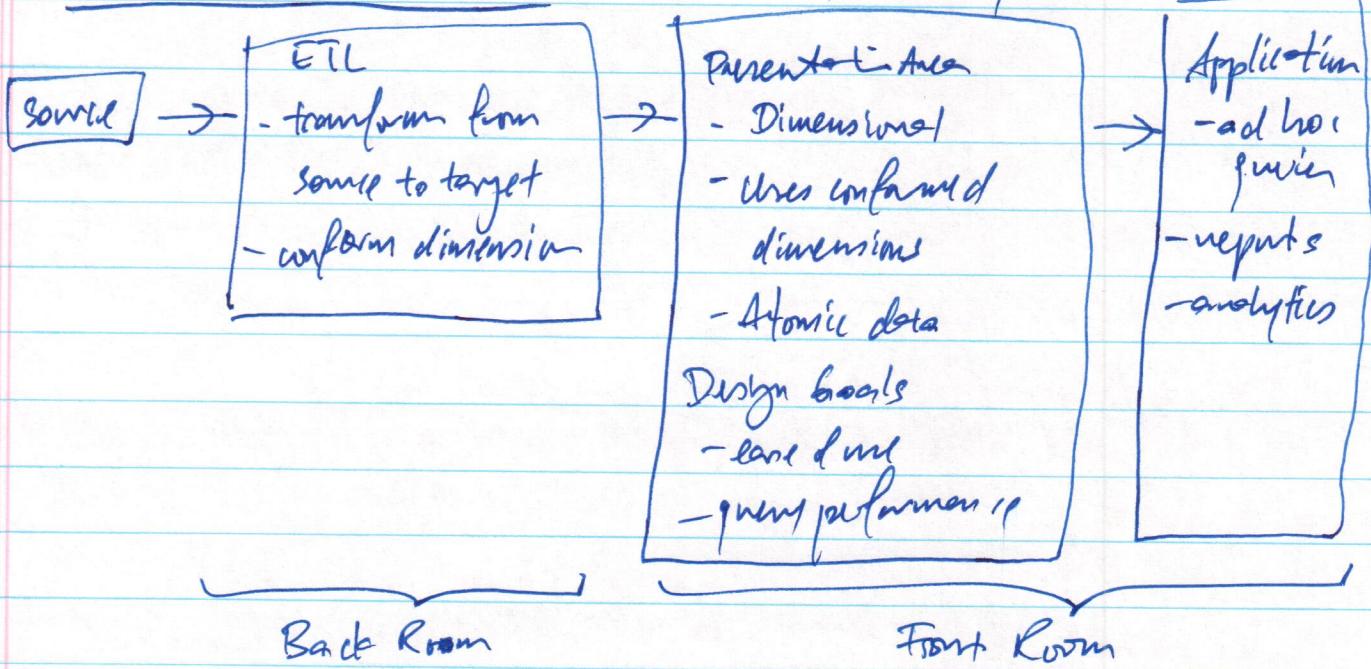
% to sql

```

Select parameter1, sum(parameter2) as x from table
group by (parameter1, parameter2, ... )
order by x desc
limit n;
  
```

Entity Relationship Diagram

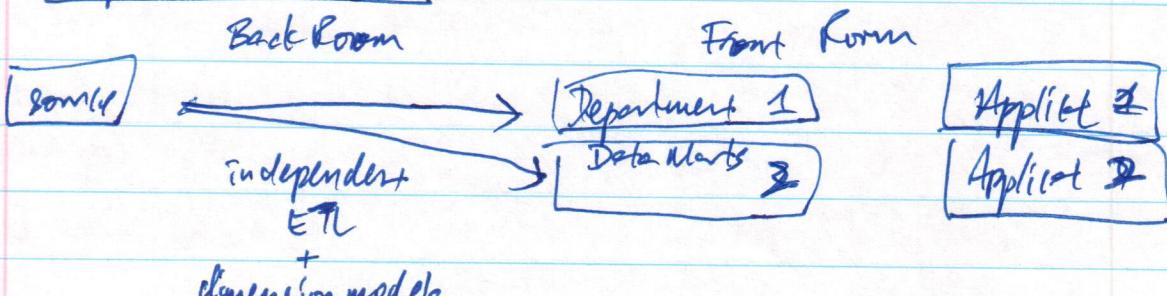
Kimball's Barn Architecture - organized by business processes



ETL

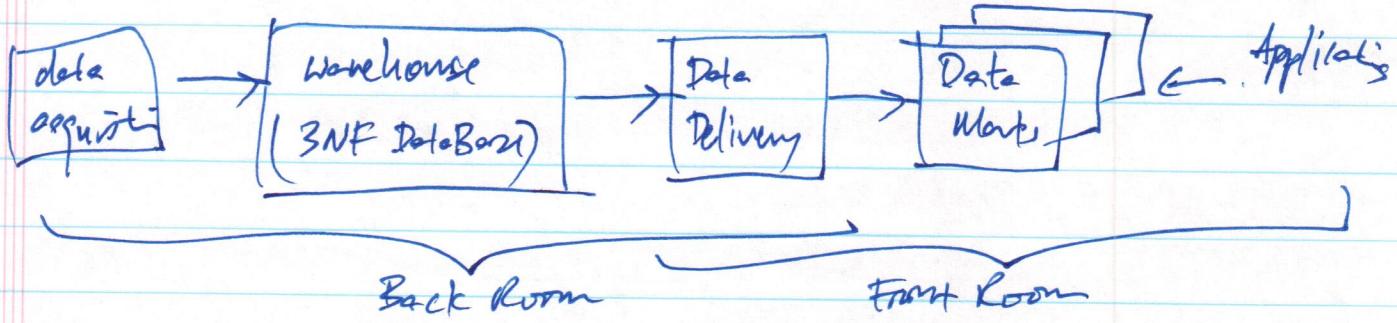
- extract - get data from sources
- transforming - integrate sources together
 - cleansing
 - diagnostic metadata
- loading - structuring + loading data into dimensional data model

Independent Data Marts



- no conformed dimensions
- uncoordinated efforts → inconsistent views

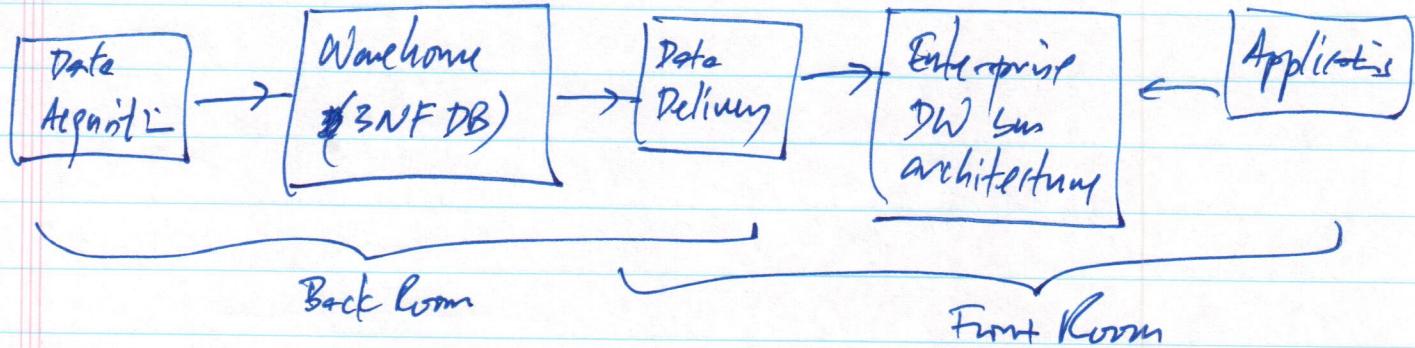
Inman's Corporate Information Factory (CIF)



• 2 ETL Process

- source → 3NF DB
- 3NF DB → data works
- 3NF DB enterprise wide data store
 - single source
- Data works dimensionally modelled (^{disimilar} ~~similar~~ to Kimball's), one aggregated.

Hybrid (Kimball Barn + Inman CIF)



OLAP cubes - aggregation of a fact metric on a number of dimensions.

Operations:

- Roll-up - sum up
- Drill-down - decompose in smaller composites down to atomic data.
- Slice - fix one dimension to single value.
- Dice - restrict some values of dimensions, same dimension cube output.

SQL - GROUP BY CUBE (x, y, z)

- one pass through fact table, aggregating all combinations of groupings.

- by x , by x,y , by x,y,z
by y , by y,z
by z , by x,z

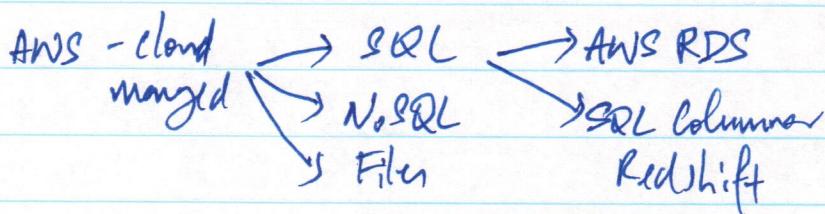
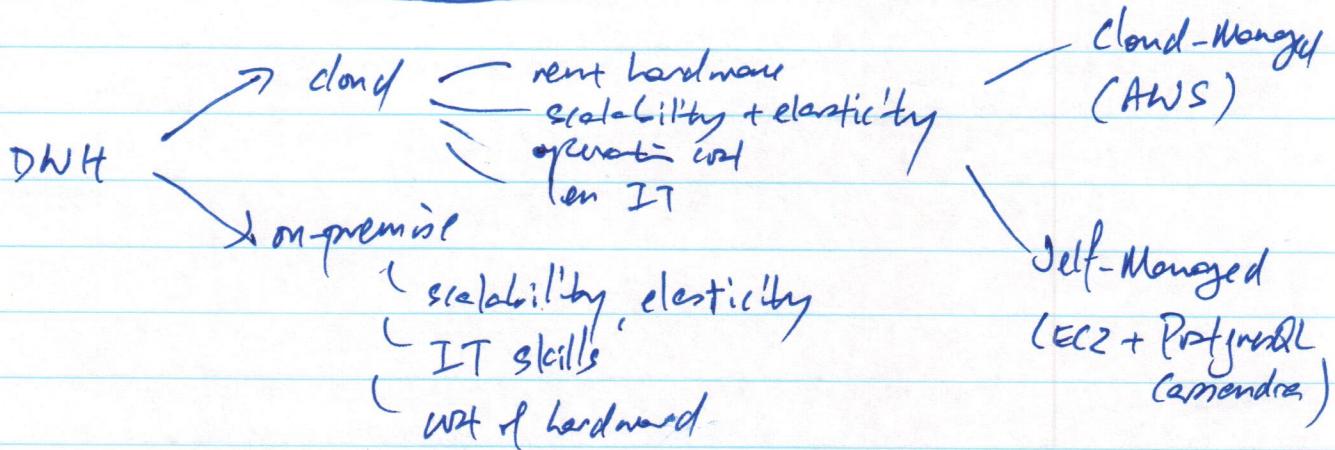
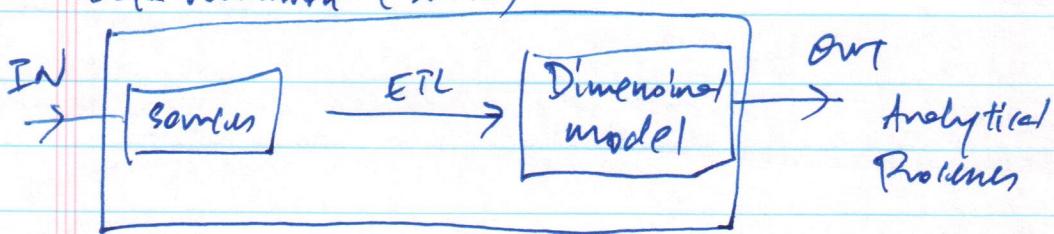
OLAP cube Technology

1. Pre-aggregate OLAP cubes, save on MOLAP (non-relational db)
2. compute OLAP cubes on the fly from existing db, (ROLAP)

Postgres - columnar storage extension (costes_fdw)
- faster than row

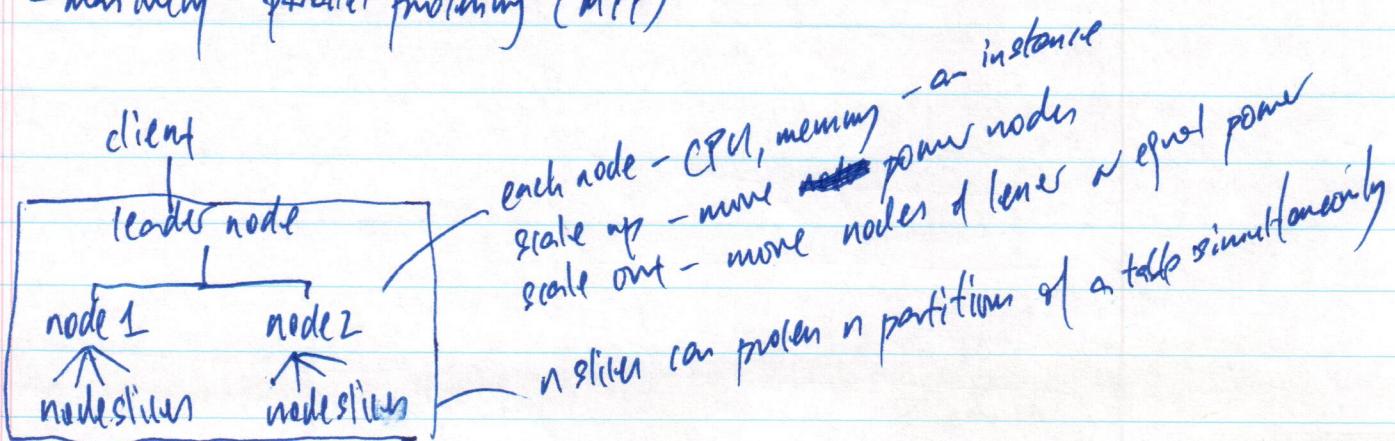
3.3 Implementing Data Warehouses on AWS

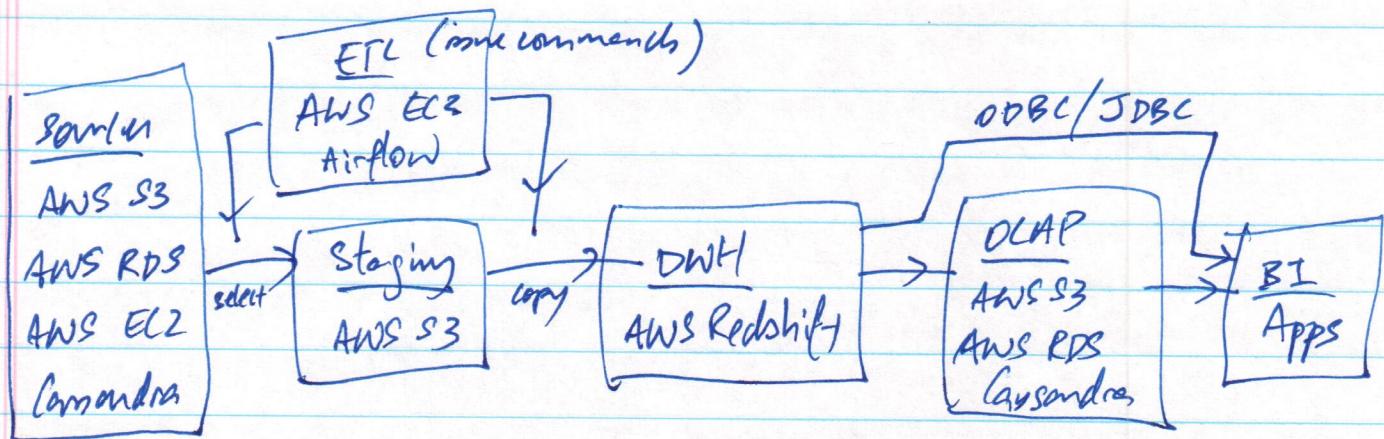
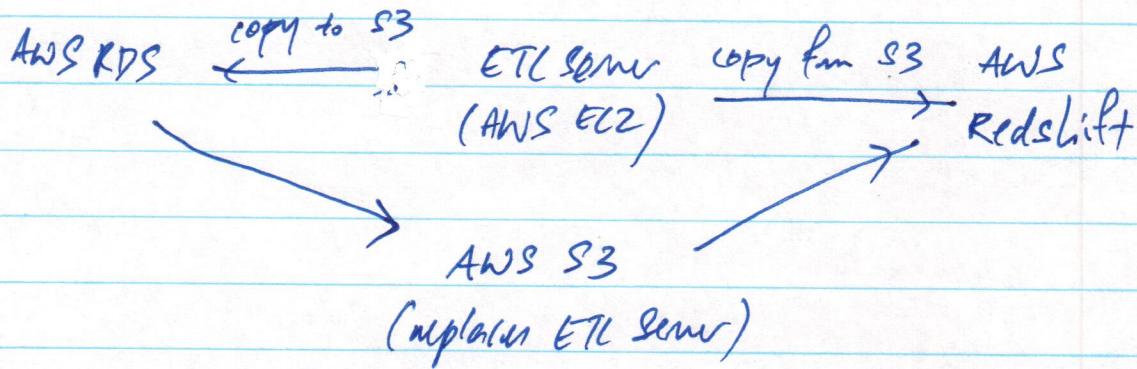
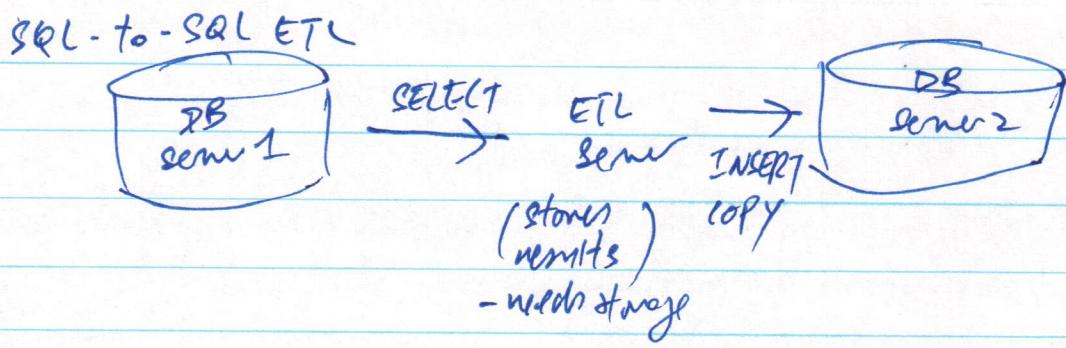
Data Warehouse (DWH)



Amazon Redshift

- column-oriented storage
- best for OLAP - summing over long history
- postgres + extensions for columnar
- run multiple queries in parallel, 1 query on $\frac{1}{n}$ CPUs
- massively parallel processing (MPP)





- Transfer data from S3 staging to Redshift with COPY
- INSERT → slow
- If file is large, split in multiple files, ingest in parallel (common prefix or use manifest file)
- compress CSV

```
COPY table FROM 's3://udacity-labs/part'
CREDENTIALS 'aws_iam_role=...'
gzip DELIMITER ';' REGION 'us-east-2';
```

```
COPY table  
FROM 's3://...'  
IAM_ROLE 'arn:aws:iam:...'  
manifest;
```

```
manifest - {  
    "entries": [  
        {"url": "s3://..."}, {"mandatory": true}  
    ]  
}
```

Ell out of Redshift

```
UNLOAD (' select * from table limit 10')  
to 's3://...'  
iam_role 'arn:aws:iam: ...';
```

Ans

- ANS

 - Redshift - IAM - 'myRedshiftRole'
 - EC2 service - ~~redshift~~ security group - 'redshift-security-group'
 - Redshift - create cluster

- Usability Data

- IAM user - airflow-redshift-wr - redshift Full Access
- credentials - .csv S3 Full Access

aws SDK = boto3

Infrastructure as Code

dwh.cfg

[AWS] (IAM user)

KEY =

SECRET =

[DWH] (create cluster)

DWH_CLUSTER_TYPE=multi-node

DWH_NUM_NODE=4

DWH_NODE_TYPE=dcl.large

DWH_IAM_ROLE_NAME=dwhRole

DWH_DB=dwh

DWH_DB_USER=dwhuser

DWH_DB_PASSWORD=Passw0rd

DWH_PORT=5439

DWH_CLUSTER_IDENTIFIER=dwhCluster

- load dwh.cfg
- create clients for EC2, S3, IAM, Redshift using boto3
- create IAM Role - iam.create_role
 - Attach Policy - iam.attach_role_policy
- create Redshift cluster - redshift.create_cluster
 - redshift.describe_cluster
- open incoming TCP port to allow cluster endpoint
- load_ext.sql
 - conn_string = 'postgresql://{}:{}@{}:{}/{}' .format(DWH_DB_USER,
DWH_DB_PASSWORD,
DWH_ENDPOINT,
DWH_PORT, DWH_DB)
↳ faster loading partitioned data (split+)
%sql \$conn_string
- ~~delete cluster~~ - redshift.delete_cluster(
 - clusterIdentifier=DWH_CLUSTER_IDENTIFIER,
skipFinalClusterSnapshot=True)

- delete role-policy - iam.detach_role_policy (RoleName = DWH_IAM_ROLE_NAME,
PolicyArn =)
delete role - iam.delete_role (RoleName = DWH_IAM_ROLE_NAME)

Optimize Table Design (over multiple CPUs)

- Distribution Styles -
even
all
auto
key

- Even - load-balancing, round-robin
 - good if a table won't be joined
 - high cost to join
- All - small tables replicated on all slices to speed up joins
 - use frequently for dimension tables
 - AKA 'broadcasting'
 - split only fact tables (even distribution)
- Auto - known design to Redshift
 - if small enough, distributed by All
 - if large, distributed by Even
- Key - rows having similar values placed in same slice.
 - skewed distribution if some values more frequent than others
 - useful if dimension table is too large for All
 - eliminates shuffling

Sorting key - especially date.
- sortkey

4. Data Cycles with Spark

4.1 Power of Spark

speed CPU → memory → SSD → network
decreasing

distributed cluster - group of nodes consists of a CPU + memory.
parallel computing - bunch of CPUs share ^{some} memory.
Partitions need in all from disk into memory.
↳ can need into chunks

Hadoop Framework

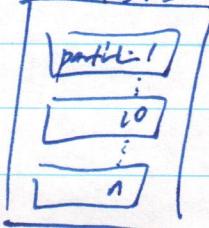
- HDFS - data storage
- map reduce - data processing
- Yarn - resource manager
- hadoop utilities

Apache Spark

Apache Storm - streaming data } in real time, low latency than
Apache Flink - streaming data } spark streaming.

MapReduce

HDFS



map
(key_i, value)

(key₁₀, value) →

(key_n, value)

shuffle

group same keys
(key₁, value)

(key₁₀, value)

(key_n, value)

reduce

(key₁, total)

(key₁₀, total)

(key_n, total)

~~ties to~~
Spark does all calculation in memory, avoiding many
data back and forth to disk.

Hadoop writes immediate calculation to disk, less efficient,
~~slow~~.

pip install numpy

Spark modes - local mode

- cluster manager - standalone
 ~ YARN
 \ Mesos

Driver (master to driver)

Don't always need Spark, if using small data that can fit into
local computer, use

- awk

- R

- Python PyData - pandas, matplotlib, numpy, scikit-learn
 \ need data in chunks.

If data in relational database, use SQL to extract, filter, aggregate
data. e.g. SQLAlchemy.

Spark is NOT a data storage system.

Data Wrangling with Spark

Scala - written, functional language

~~Applies to~~ API - Python, JAVA, R → PySpark

Functional programming - perfect for distributed computing
Spark does not change input data, makes a copy, ie immutable.
Lazy evaluation - build step-by-step description of what functions + data needed

→ Directed Acyclic Graph (DAG)
- prognostic running func to get data till ~~func~~ end.
full ~~map~~ processing using 'collect()'

```
import findspark
findspark.init('spark-2.3.2-bin-hadoop2.7')
```

```
import pyspark
sc = pyspark.SparkContext(appName='tasty')
log_of_songs = []
distributed_log_songs = sc.parallelize(log_of_songs)
distributed_log_songs.map(lambda song: song.lower()).collect()
```

Data Formats

- CSV, JSON, HTML, XML,
- HDFS - Hadoop Distributed File System/Storage
 - split into 64 or 128 MB blocks
 - replicate blocks into cluster
 - fault tolerance + digestible chunks
 - AWS S3 - store raw data

```
from pyspark import SparkContext, SparkConf
conf_name = SparkConf().setAppName('name').setMaster('IP Address')
sc = SparkContext(conf=conf_name)
```

```
from pyspark.sql import SparkSession  
spark = SparkSession.builder.appName('appName')  
    .config('config option', 'config value')  
    .getOrCreate()
```

```
spark.SparkContext.getConf().getAll()
```

```
user_log=spark.read.json(path)
```

```
spark.read.csv(path, header=True)
```

Imperative Programming

- ~ spark dataframes + Python
- ~ how?

Declarative Programming

- ~ SQL
- ~ what?
- ~ abstractive layer on top of imperative programming.

Spark DataFrames

- select()
- filter()
- , where()
- , groupBy()
- , sort()
- , dropDuplicates()
- , withColumn()

Spark SQL

- ~ user defined functions
- ~ window functions - partitionBy, rangeBetween or rowBetween

```
user_log.createOrReplaceTempView('user_log_table')
```

```
spark.sql('SELECT * FROM user_log_table').show()
```

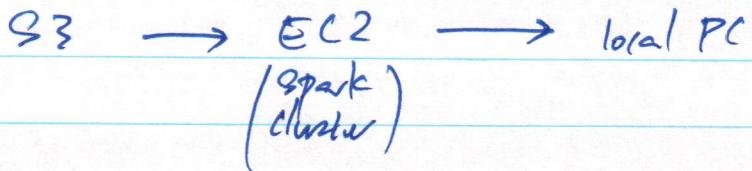
```
spark.udf.register('funct', lambda x: x.lower())
```

spark query optimizer - catalyst - translate to DAG from either
pyspark or SQL

DATA execution plan on RDD (resilient distributed data set)

Spark Cluster Managers

- Standalone
- MESOS } sharing a rm team
- YARN }



Install Spark individually to each machine, manual, painful,
prone to user error.

- AWS EMR (Elastic Map Reduce) does all this automatically

HDFS uses MapReduce as a resource manager for distributed.

Spark lacks this; but is faster.

* Install Spark on HDFS.

Spark Cluster

1. EC2 - key pair → .pem

2. EMR - create cluster

emr-5.20.0 + spark 2.4.0 + Hadoop 2.8.5
m5.xlarge

spark-submit --master yarn ./lower_songs.py

AWS S3 - object storage system, stores all kinds of format.

HDFS - distributed file system, guarantees fault tolerance via duplicates,
only format avro and parquet

4.4 Debugging + Optimizing

_corrupt_record - populates if input data is bad.

log-data.where(log-data['_corrupt_record'].isNotNull().collect())

Accumulators - for debugging.

incorrect_records = SparkContext.accumulator(0, 0)

incorrect_records.value

def add_incorrect_record():

global incorrect_records

incorrect_records += 1

Spark WebUI - DAbs
Stages
Tasks

SSH - port 22

HTML - port 70

Jupyter notebook - port 8888

Spark Jobs - port 4040

WebUI - port 8080

Transformation faults - select, filter → don't show up in tasks

Action function - write, head → show up in tasks

sparkContext.setLogLevel("ERROR")
("INFO")

- Data skew
- Pareto Principle
 - identify by sampling + summarizing (min, max, identify outliers, many long)
 - change workload division by timestamp or user or ...
 - partition data
 - ↓
 - df.repartition(
number
of rows)
 - make composite keys

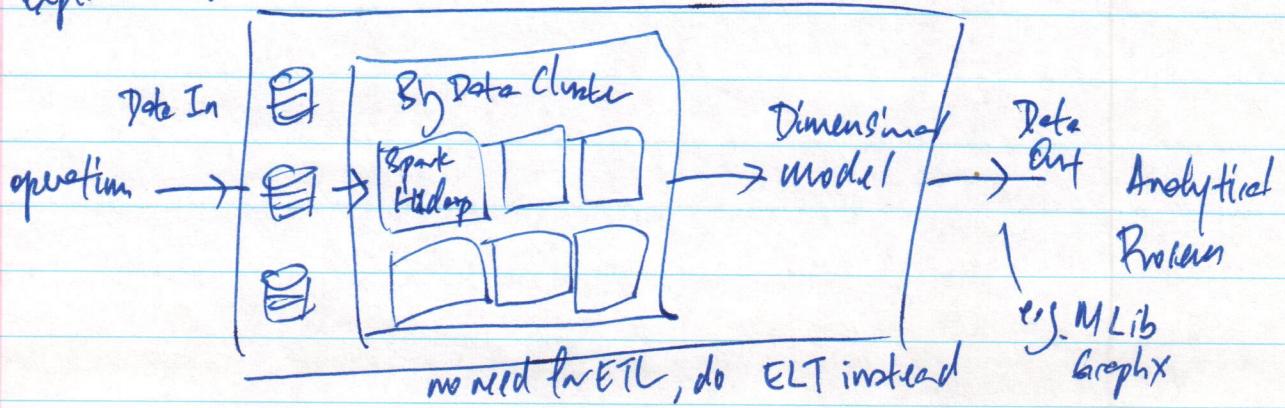
Out of memory errors, increased number of partitions
Inefficient growth.

4.5 Data Lakes

- unstructured data (text, xml, json, ...)
- volume (IoT, social)
- big data technologies (HDFS, Spark)
- new data analysis (deep learning, ...)

- harder to define into facts / dimension tables / json schema-on-read - no creation of tables
- cheaper with HDFS than mPP

freedom to expand beyond existing tables / schema, ad-hoc data exploration.



Schema - On - Read

- skip datalane with hadoop / Spark , no insertions
- read file immediately

```
spark.read.csv("... .csv",
    inferSchema = True,
    header = True,
    sep = ";",
    mode = "DROPMALFORMED")
```

Un- / Semi- / Structured Input

- spark.read.text ("... .pz") } text based
- spark.read.csv ("... .csv") }
- binary - Avro , Parquet
- compressed
- negrid
- write to local , HDFS , S3
- read/write to databases
 - ~ SQL through JDBC
 - ~ NoSQL - mongoDB , cassandra ,

AWS

Storage

HDFS

S3

S3

Presto

Spark

Spark

Dynamodb

ANS-Managed

EMR (HDFS + Spark)

TMR (Spark only)

Athena

Vendor-Managed

EC2 + vendor soluti.

EC2 + "

Services + "

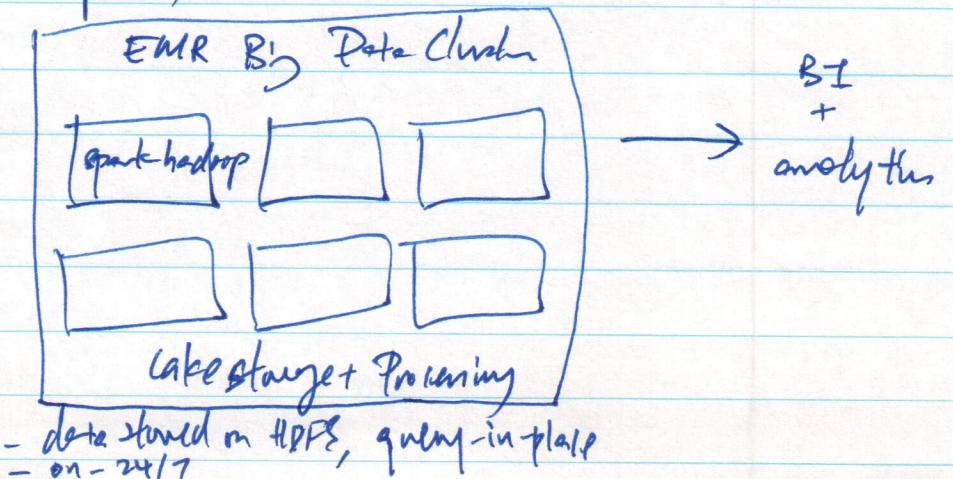
	Data Warehouse	Data Lake
Data Form	tabular	all formats
Data value	high value	high, medium
Ingestion	ETL	ELT
Data model	Star, snowflake schema w/ conformed dimensions or delta-marks and OLAP cubes	Star, snowflake, OLTP and-hoc representation
Scheme	schema-on-write	schema-on-read
Technology	Expensive MPP database	Commodity hardware w/ parallelism
Data Quality	high with effort for consistency + clear rules for accessibility	Mixed - raw data, some transformed to higher quality
User	Business analyst	Data scientist, Business analyst, ML engineer
Analytics	Report + BI	Machine learning, graph analytics, data exploration.

AWS EMR (HDFS + Spark)

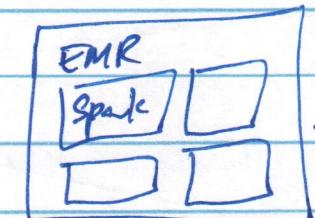
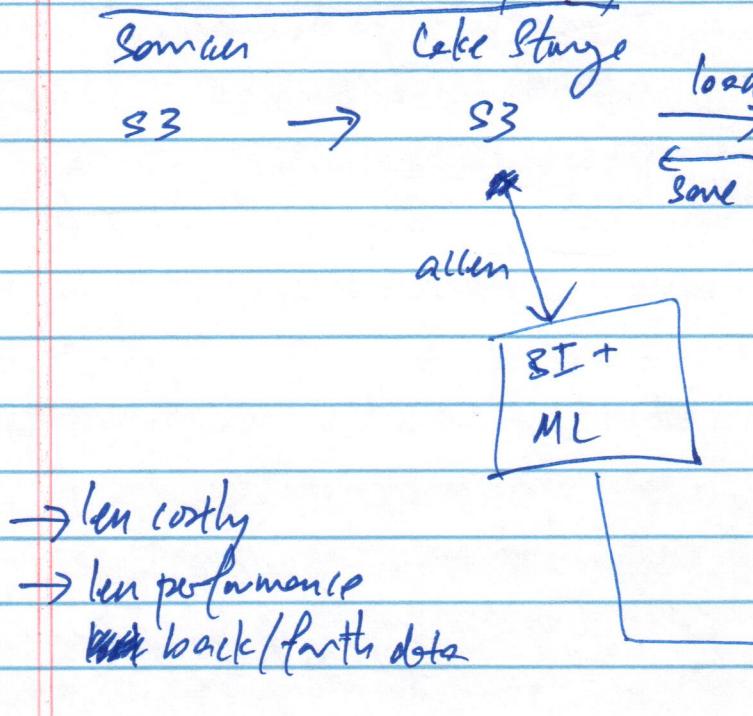
Samuel

ss

Import



AWS EMR (S3 + Spark)



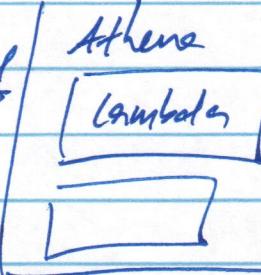
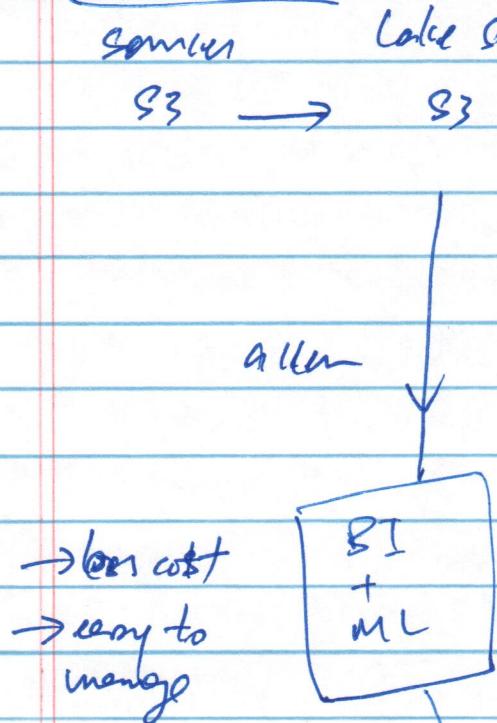
EMR Big Data Cluster
(Lake Processing Only)

- no HDFS
- all data stored in S3
- problem in EMR

- results saved back to S3
- can turn off (no need 24/7)

all the

AWS Athena



Lambda services

(Lake Processing Only)

- stored on S3
- load + process data on serverless lambda execution (pay by execution time, not up time)

all the

```
import pyspark.sql.functions as F  
df = df.withColumn('payment_date', F.to_timestamp('payment_date'))
```

```
from pyspark.sql.types import StructType as R  
StructField as Fld  
DoubleType as Dbl  
StringType as Str  
IntegerType as Int  
DateType as Date  
paymentSchema = R([Fld('payment_date', Int()),  
"",""],  
])
```

```
spark.read.csv("s3://... ", sep=";", schema=paymentSchema,  
header=True)
```

AWS Glue crawler - scheduled ETL

- point to S3 raw data
- schema-on-read

AWS EMR + notebook

spark cluster on AWS EC2

5. Data Pipeline w/ Airflow

5.1 Data Pipeline

- series of steps in which data is processed.

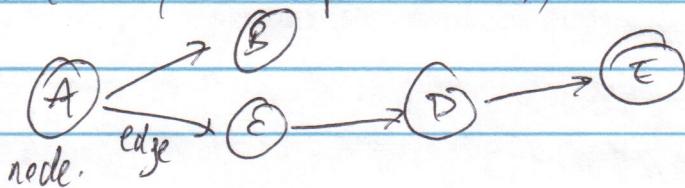
Kafka - open source stream processing software platform for LinkedIn.

Data validation - process of ensuring data is present, correct + meaningful.

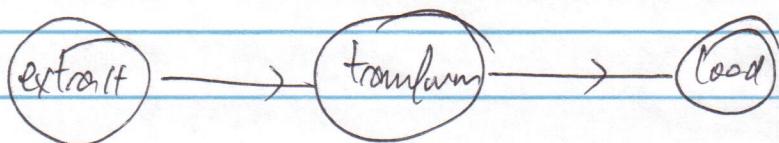
- automated checks.

- valid number of rows in Redshift epochs in S3

Directed Acyclic Graphs (DAGs)

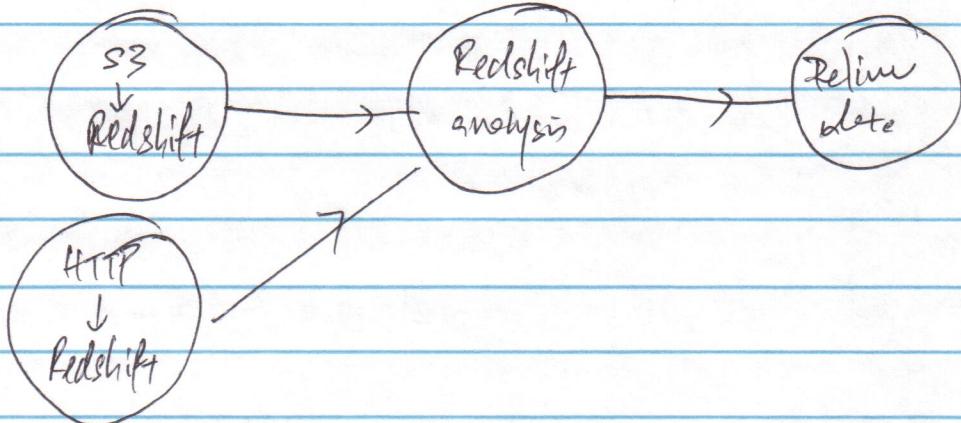


- do not contain cycles
- have direction



ETL - each step is node

- dependent on prior steps are directed edges



Airflow - open source tool to structure data pipeline as DAGs

- Airbyte in 2015

- DAGs written in Python, run on schedule on/fair cluster
- integrate with external tools Redshift, spark, Hadoop etc easily.

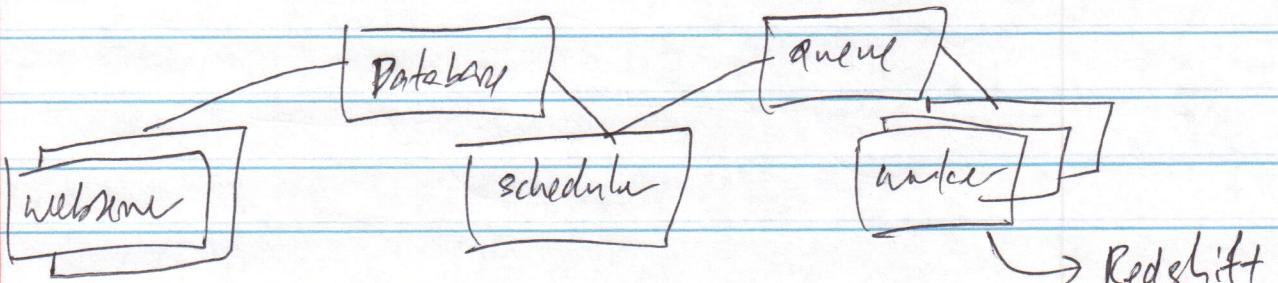
- scheduler executes tasks in array of workers

- web-based UI to visualize + interact with data pipeline

```
[ ] import logging  
from airflow import DAG  
from airflow.operators.python_operator import PythonOperator
```

```
def helloworld():  
    logging.info("helloworld")
```

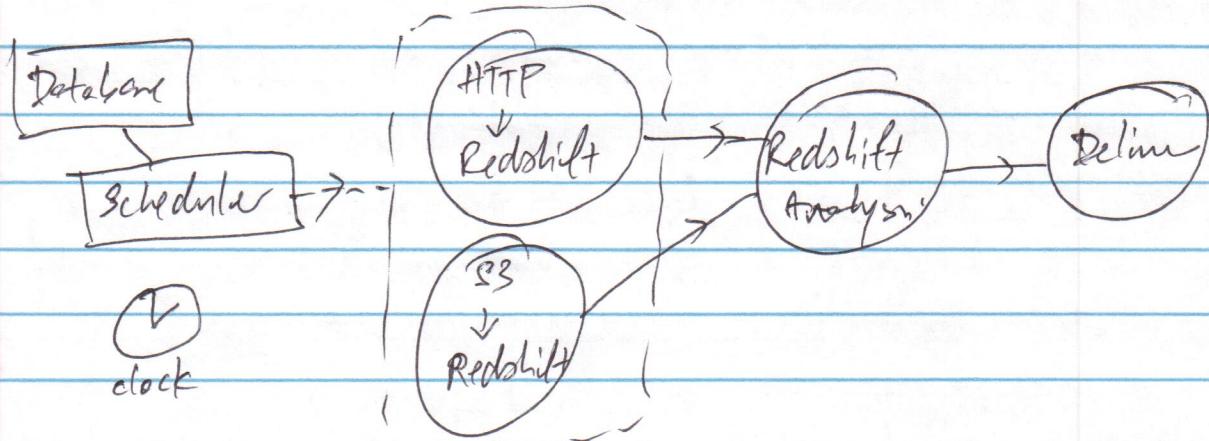
```
dag = DAG('lesson1-exercise1',  
          start_date=datetime.datetime.now())  
greet_task = PythonOperator(task_id="hello-world-task",  
                           python_callable=helloworld,  
                           dag=dag)
```



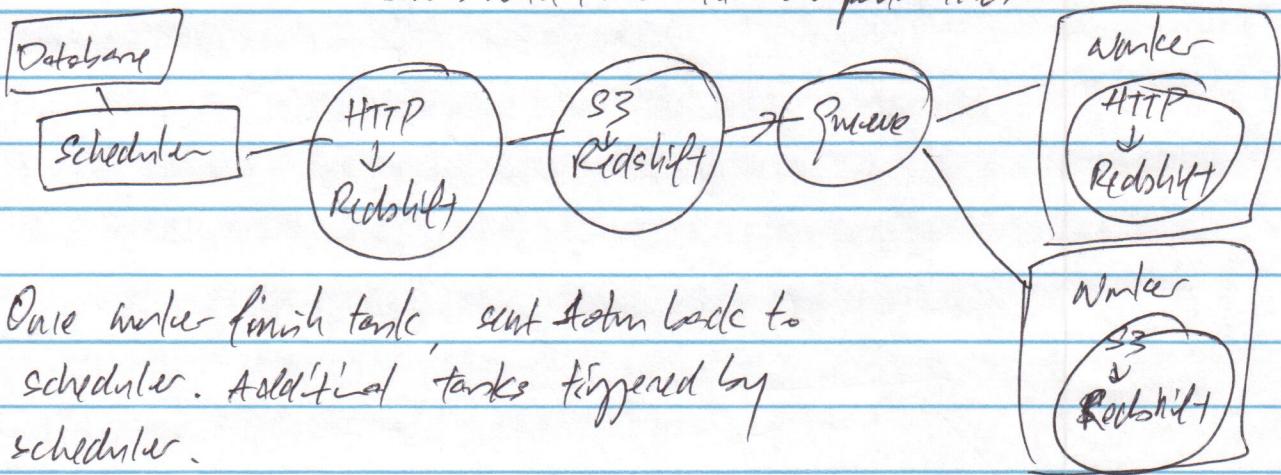
- database - credentials, connections, history
- webserver - control dashboard
- Redshift, Spark, etc...

- coordinate movement of data between data storage and data processing tools.
- don't pass data in memory, not data processing framework.
- limit to 10 workers in 1 machine.

Airflow Scheduler starts DAGs based on time & external triggers.



starts with those with no dependencies.



One worker finish task, sent status back to scheduler. Additional tasks triggered by scheduler.

Creating a DAG

```
from airflow import DAG, from airflow.operators.python_operator  
import PythonOperator
```

```
DAG → divvy_dag = DAG('divvy',  
                      description = "analyze data",  
                      start_date = datetime(2019, 2, 4),  
                      schedule_interval = '@daily')
```

```
Operator → task = PythonOperator(task_id = 'hello_world',  
                                 python_callable = helloworld,  
                                 dag = divy_dag)
```

Schedulers are optional - (Q) only

(Q) hourly

(Q) daily

(Q) weekly

(Q) monthly

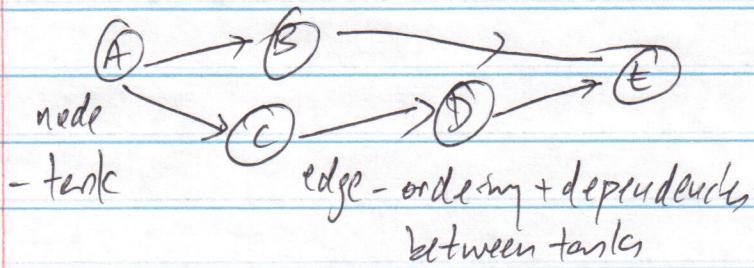
(Q) yearly

None

CRON strings

- AIRFLOW OPERATORS
- Python Operator
 - Postgres "
 - RedshiftToS3 "
 - S3ToRedshift "
 - ...

- Bash Operator
- SimpleHttp "
- Sensor
- ...



$a >> b$, a before b
 $a << b$, a after b

```
hello_world_task = PythonOperator(....)
goodbye_world_task = PythonOperator(....)
hello_world_task >> goodbye_world_task
or      "      " .set_downstream( " " )
```

Connections via hooks - reusable interface to external systems and databases.

e.g. db_hook = PostgresHook('demo')

df = db_hook.get_pandas_df('SELECT * FROM rides')

HTTP Hook

Postgres Hook (Redshift)

MySQL Hook

Slack Hook

Presto Hook

```
hook = S3 Hook (aws_conn_id='aws_credentials')
```

```
bucket = Variable.get('s3_bucket')
```

```
prefix = Variable.get('s3_prefix')
```

```
def hello_date(*args, **kwargs):
    print(f"Hello {kwargs['execution_date'])")
```

```
dinvys_dag = DAG(...)
```

```
task = Python Operator(task_id='hello_date',
                       python_callable=hello_date,
                       provide_context=True,
                       dag=dinvys_dag)
```

```
def load_date_to_redshift( ):
```

```
aws_hooks = AwsHook('aws_credentials')
```

```
credentials = aws_hooks.get_credentials()
```

```
redshift_hook = RedshiftHook('redshift')
```

```
redshift_hook.run(sql_statements=COPY_ALL_TRIPS_SQL.format(credentials.access_key, credentials.secret_key))
```

```
create_table = PostgresOperator(task_id='create_table',
```

```
dag=dag, postgres_conn_id='redshift',
```

```
sql=sql_statements.CREATE
```

copy_task = PythonOperator(task_id='load_from_s3_redshift',
dag=dag,
python_callable=load_data_to_redshift)
create_table >> copy_task

3.2 Data Quality

~~Data lineage~~ - discrete steps involved in creation, movement, calculation of dataset.

- instilled confidence
 - define metrics
 - helps in debugging
 - Airflow DAGs are natural representation for transformation + movement of data
- e.g. create_table → load_to_redshift → calculate

schedules - based on size of data

- frequency of data
- related datasets

Backfill - automatically analyze.

End-date - end of life

DAG - max active runs

Data partitioning - process of isolating data to be analyzed by one or more attributes (time, logical type, data size)
- for speed, more reliable pipelines.

plugins/operators/`__init__.py`
`/gs3_to_redshift.py`

DAG tasks - atomics, have single purpose.

- maximize parallelism
- make failure states obvious.

SubDAGs - commonly repeated series of ~~task~~s within DAGs

- reduce code, debug, easier to understand.
- limits visibility in Airflow UI, ~~but~~ obfuscate

SLA - email alerts

- publish metrics to statsd