

Udacity: Data Engineering with Azure

Data Modeling

1. Introduction to Data Modeling

Data modeling - Relational Database Management System (RDBMS)

↳ PostgreSQL

↳ MySQL

- NoSQL

↳ Cassandra

↳ MongoDB

↳ denormalized
queries first - 1 table per query
can't join tables

OLAP vs OLTP

Normalized vs Denormalized tables

↳ 3NF

↳ fast read
slow write

Fact tables, Dimension tables

↳ star

↳ snowflake

Cloud Data Warehouses with Azure

1. Introduction to Cloud Data Warehouses with Azure

Data Sources → Data Engineering → Data Consumers

- CSV - video, audio

- data pipelines

- analytics

- JSON

- ingest

- visualization

- database

- transform

- data science

- logs, systemstate,

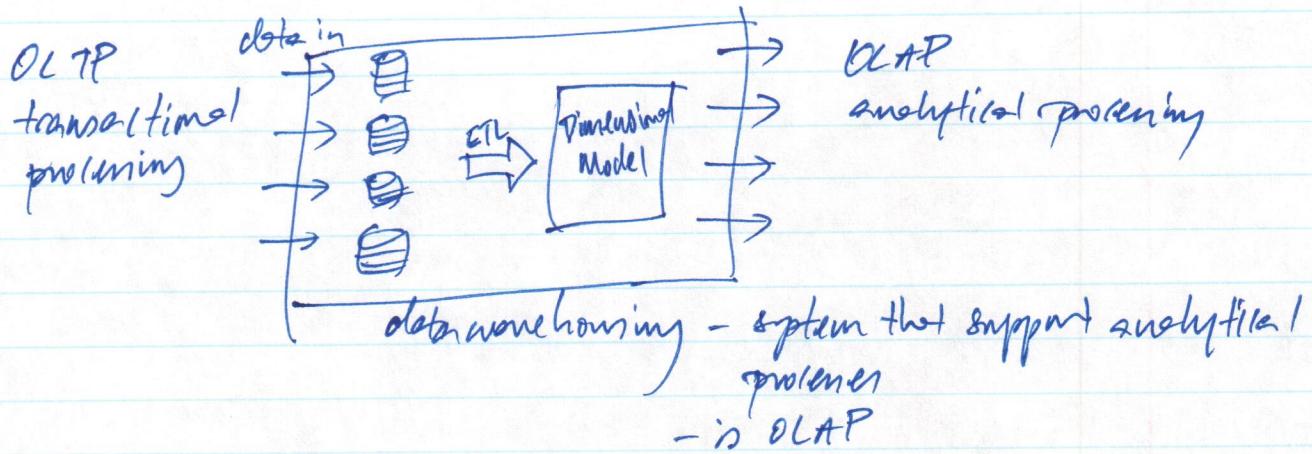
- store

2. Introduction to Data Warehouses

In Business:

Operational Processes - customers, inventory, delivery
vs.

Analytical Processes - HR, marketing, management.



If same database for OLAP and OLTP, database scheme too complex for business analysts
- analytic queries require too many joins
- ok if database is small.

Data Warehouse:

- Kimball

- is a copy of transaction data specifically structured for query + analysis

- Inmon

- is subject-oriented integrated, nonvolatile, time-variant collection of data in support of management's decision.

- Rainbird

- is system that gathers + consolidates data periodically from source into a dimension or normalized data store. Keeps years of history. Query for business intelligence or other analytical activities. Updated in batches.

Dimensional Modeling:

3NF - expensive joins

- hard to explain to business users

Star schema - joins with dimensions only

- good for OLAP, not OLTP

- fast query

Facts table

- quantifiable business event or metric

Dimension table

- context of business event - e.g. who, what, why, where, etc.

3NF to Star schema - ETL

* extract - join tables

- change types

- add new columns

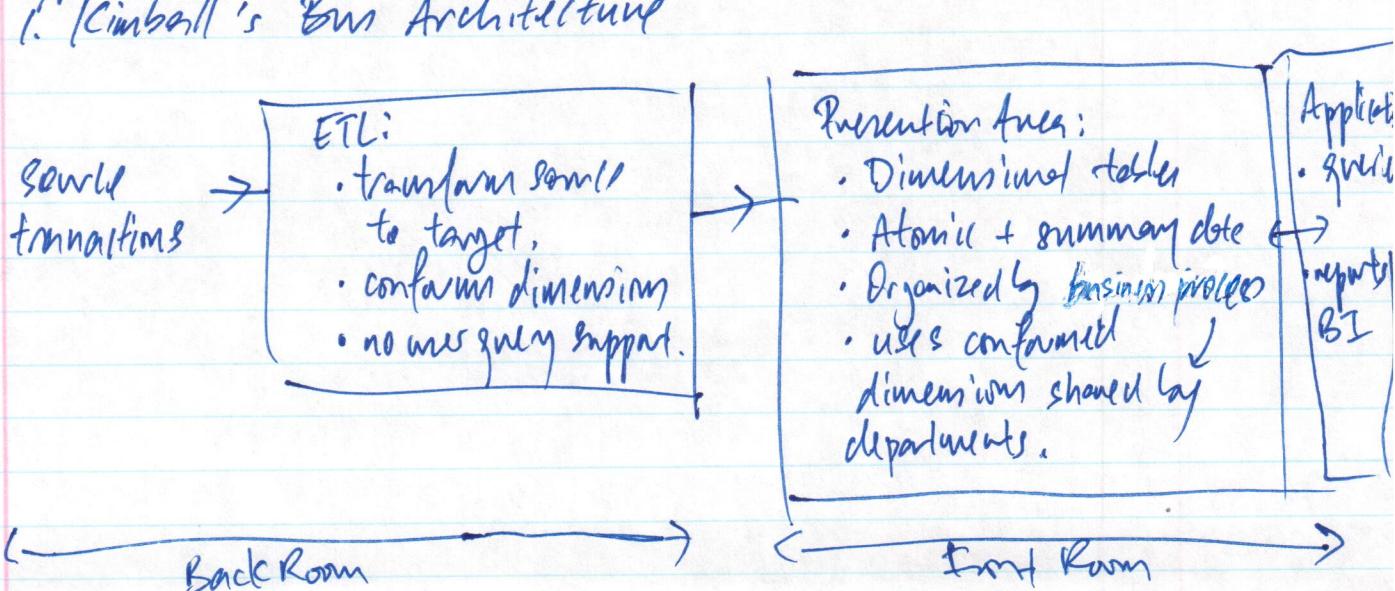
* load

- insert into facts, dimension tables.

} transform

Types data warehouses:

1. Kimball's Bus Architecture



Extract

- transfer data to warehouse.

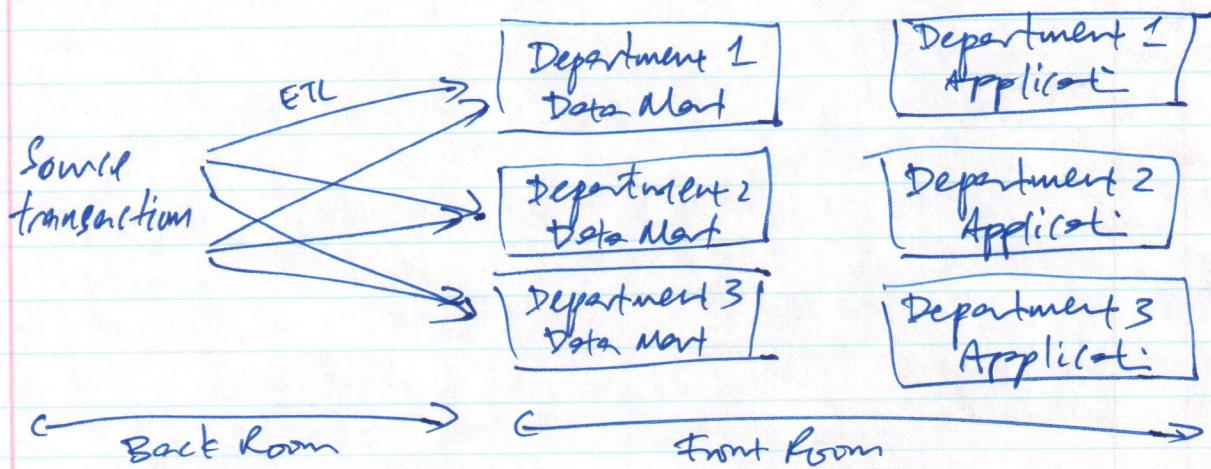
Transforming

- integrate many sources together.
- cleansing - duplicates, inconsistencies, missing values.
- produce diagnostic metadata.

Loading

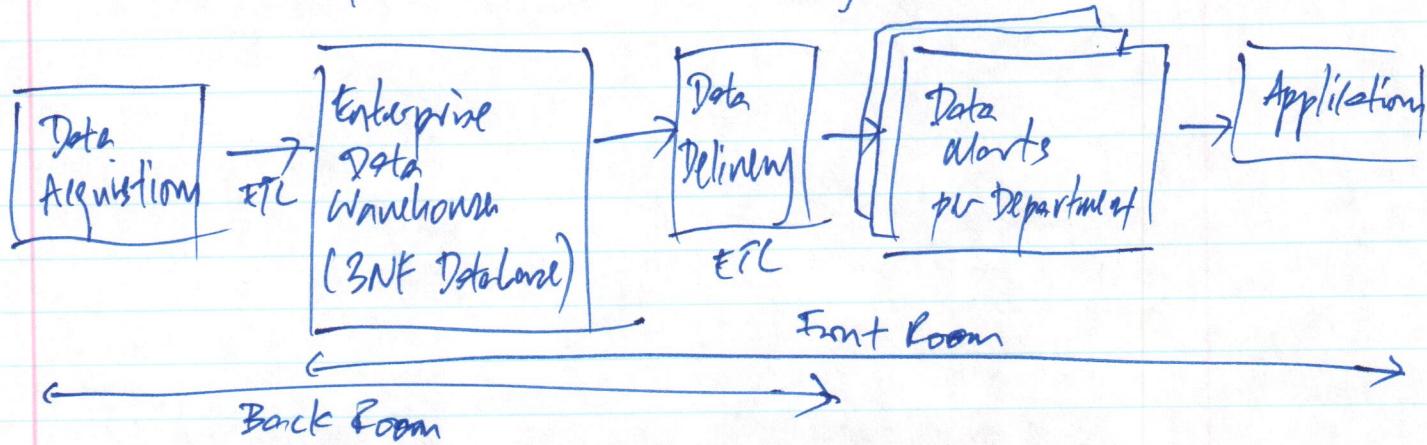
- structuring + loading data into dimensional data model!

2. Independent Data Marts



- separate ETL processes + dimensional models called data marts.
- different fact tables for same events, no conformed dimensions.
- generally discouraged.
- can lead to inconsistent views.

3. Inmon's Corporate Information Factory (CIF)



- 2 ETL processes

1. Source → 3NF database

2. 3NF database → departmental data marts

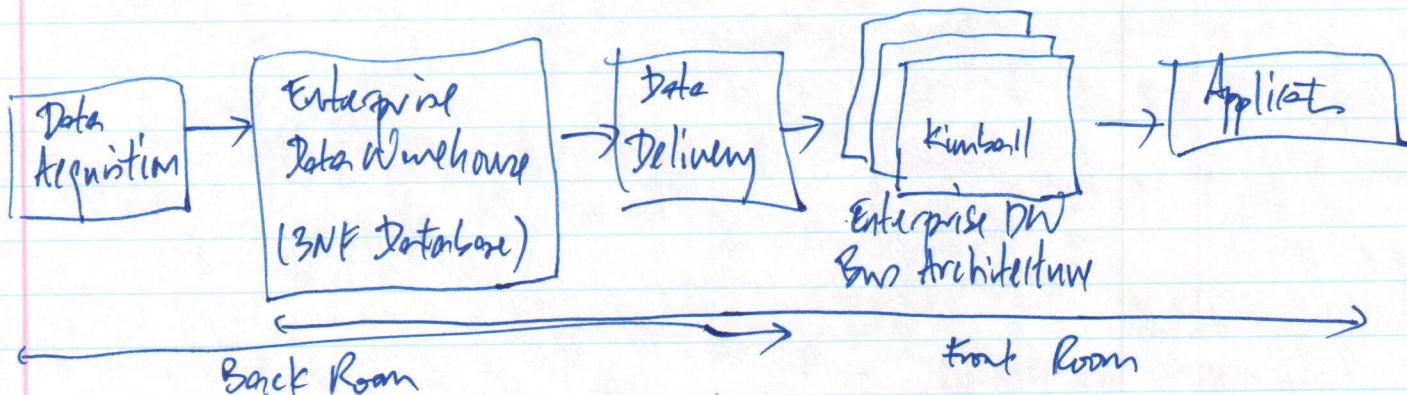
- 3NF database acts as enterprise-wide data store

- single integrated source of truth for data marts

- could be accessed by end-users if needed

- data marts are dimensionally modeled, mostly aggregated, unlike Kimball

4. Hybrid Kimball Bus + Inmon CIF



- Maintains 3NF Database for ground truth.

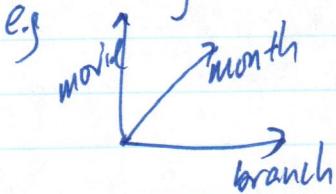
- Removed data mart

- Replaced with enterprise data warehouse (Kimball - confirmed dimension).

OLAP Cubes

in star schema

- After creating facts + dimension tables, create OLAP cubes.



- OLAP cube is aggregate of number of dimensions
- rollup, drill-down, slice + dice.

- Rollup - sum up columns.

- Drill down - decompose into smaller grains.

- Slice - reduce N dimensions to $N-1$ where N^{th} dim is constant value.

- Dice - compute sub-cube - restrict range of values in columns, while retain same N dimension.

- GROUP BY CUBE (x, y, z) will make one pass through fact table, and aggregate all possible combinations of groupings of lengths $\emptyset, 1, 3$

e.g. total revenue, revenue by x , revenue by x, y , revenue by x, y, z
" by y ; " by z ; " by z ;

- Build OLAP cubes

1. Pre-aggregate OLAP cubes, some on non-relational database (MOLAP)
old approach.
2. Compute cube on the fly, ROLAP

3. ELT and Data Warehouse Technology in Cloud

Benefits of cloud DW over on-premise,

1. Scalability

- expand with more data.

2. Flexibility

- add or remove services from pipeline as needed.

- change to accommodate business needs changes.

3. Cost Shifting

- perform costly + time consuming transformation of data as last part of process for 'just-in-time' business prioritization.

ETL vs ELT

- ETL - transformation happens on intermediate server

- ELT - " " " destination server via stepping tables.

- better scalability

" flexibility

" cost shifting

" performance on large datasets

more flexibility for NoSQL datasets

Relational Databases:

- oracle
- Microsoft SQL Server
- PostgreSQL
- MySQL
- MariaDB

Microsoft Azure:

- Azure SQL Database
- " " Database for MySQL
- " " " MariaDB
- " " " PostgreSQL

GCP -

- cloud SQL (MySQL, PostgreSQL, MS SQL server)

AWS

- Amazon RDS (MySQL, PostgreSQL, MariaDB, Oracle, MS SQL)

NoSQL database:

- key value
- document
- column oriented
- graph
- time series

Azure

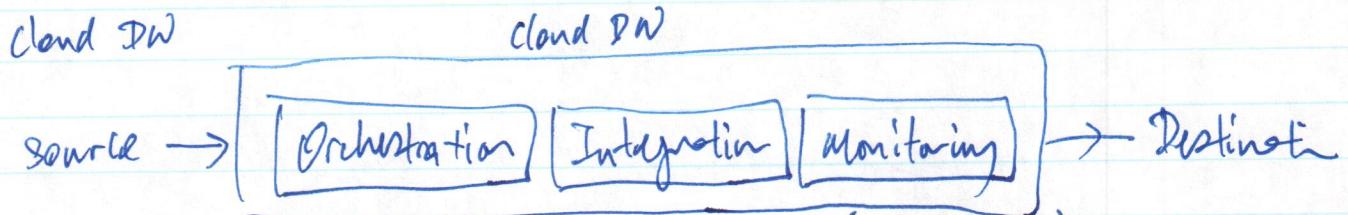
- CosmosDB
- Gremlin - graph database
- MongoDB - document
- Cassandra - column oriented

GCP

- Bigtable - column oriented
- Firestore - document
- MongoDB atlas - document

AWS

- DynamoDB - key value
- Document DB - document
- Keyspace - column oriented
- Neptune - graph
- Time stream - time series



ETL/ELT

- Azure Data Factory
- AWS Glue
- GCP Dataflow

Streaming Data (IoT, web)

- Azure - Streaming Analytics
- AWS - Kinesis
- GCP - DataFlow

3rd party companies providing ETL/ELT

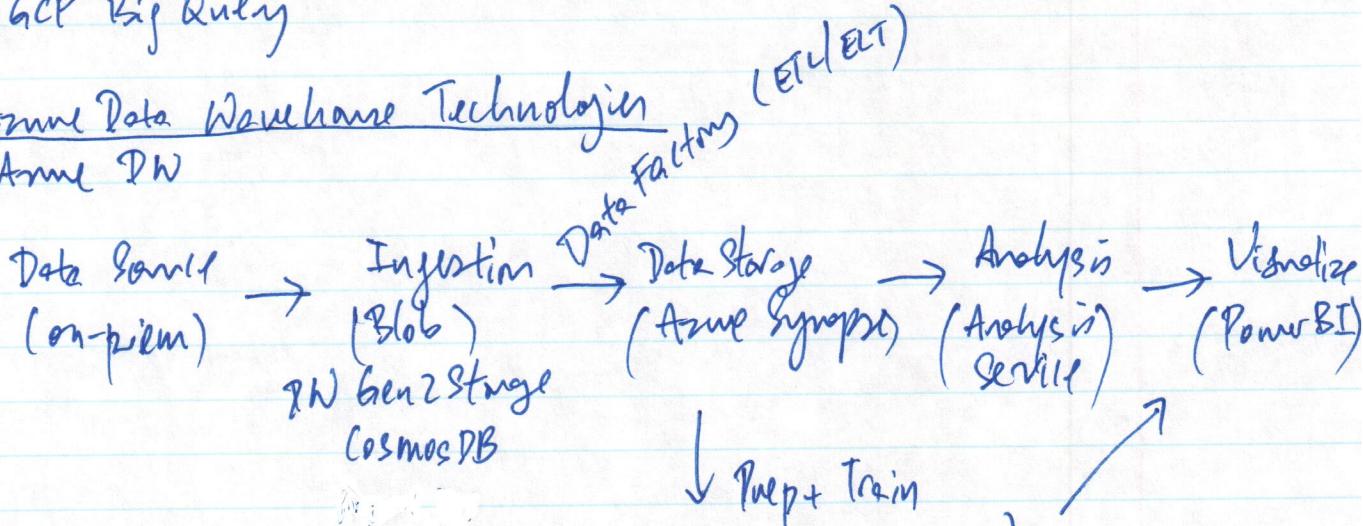
- Informatica
- Talend
- xplenty
- Matillion

Cloud DW Solution

- Azure Synapse
- AWS Redshift
- GCP BigQuery

4. Azure Data Warehouse Technologies

Azure DW



DS: Data Source

Azure Portal → Synapse Analytics workspace → Open Synapse Studio



Data - tabular databases

Develop - Add notebooks to query data - e.g. SQL scripts.

Integrate - add pipelines from data sources

Monitor - look at runs for each of above.

Manage - manage components Synapse

Cloud Data Storage:

- Azure Data Warehouse Gen 2 - traditional DW architecture
- Azure Dedicated SQL Pools - relational data storage
- Blob storage - file based storage
- CosmosDB for NoSQL - column-oriented, or document db.

ETL/ELT Pipelines:

- Azure Data Factory - data integration, data flows
- Azure Databricks - use Spark to create ETL pipelines / SQL
- Azure Polybase - use TSQL to query blob storage to support ETL

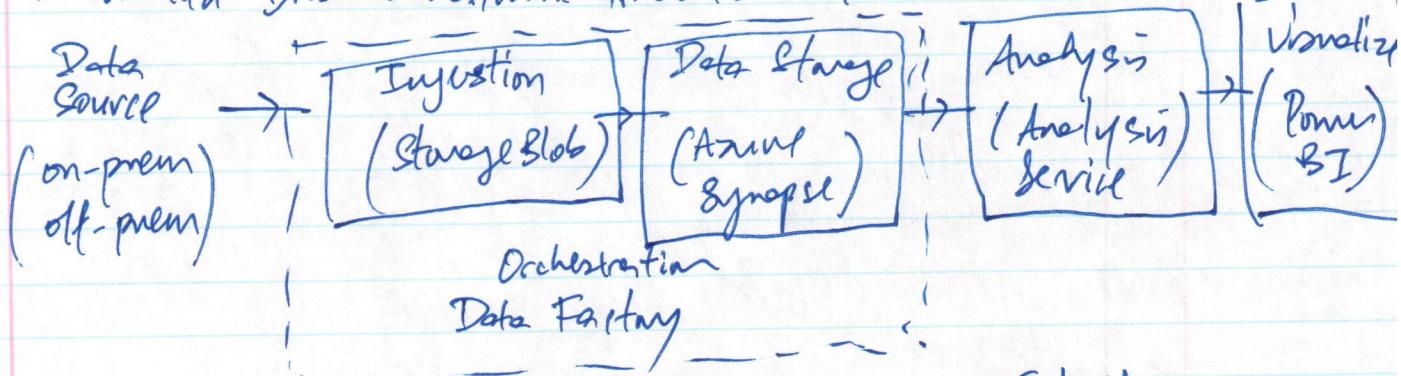
python

Save

SQL

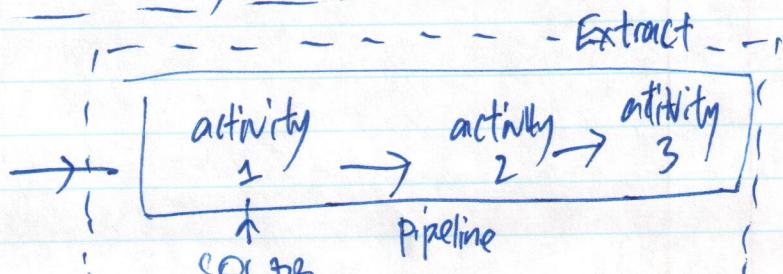
Implementing Data Warehouses in Azure

Automated Data Warehouse Architecture



Data Ingestion:

Source data → trigger → Linked Service



Create Linked Services:

Azure Portal → Azure Synapse Workspace → Open Synapse Studio

↳ managed - external connection

↳ linked services

↳ + New

↳ select type

↳ Name

↳ From Azure Subscription

↳ manually

username / password

Test Connection ↙

Import - from source to destination

- Built in copy task

- Run Once / Schedule / Tumbling Window

- Source type

Connection

↳ source folder

↳ Preview

↳ Target type

Connection

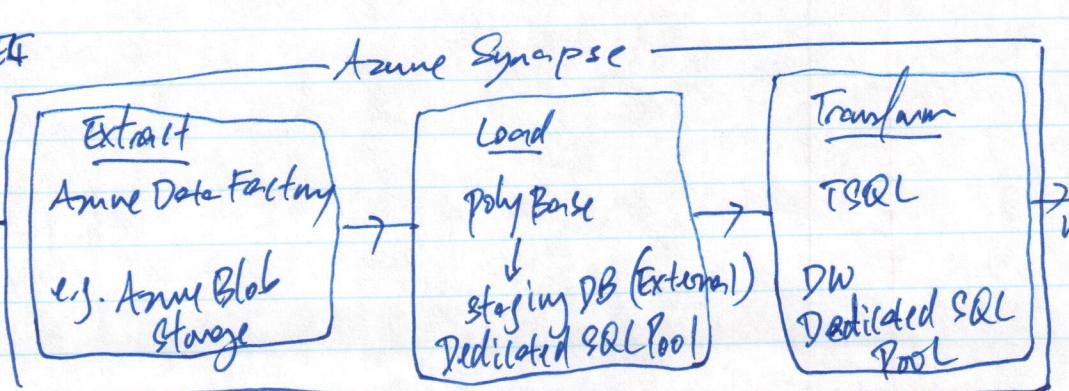
↳ destination

file type

→ View in Data

SQL to SQL ETL

Source data



Create Staging Tables:

Synapse Studio → manage → Analytics pool → SQL pools
dedicated

Data → Workspace → Tables

→ External tables (staging tables)

→ Linked → Azure Data Lake Storage Gen 2

↳ Blob Storage

↳ New SQL Script

↳ Select TOP 100 rows

Create External Table

Bulk load

File type

Select SQL pool

Select database

External table name

Create external table

↳ Automatic

↳ Use SQL script

Open script

Edit

Run

view to check

DROP table

CREATE table

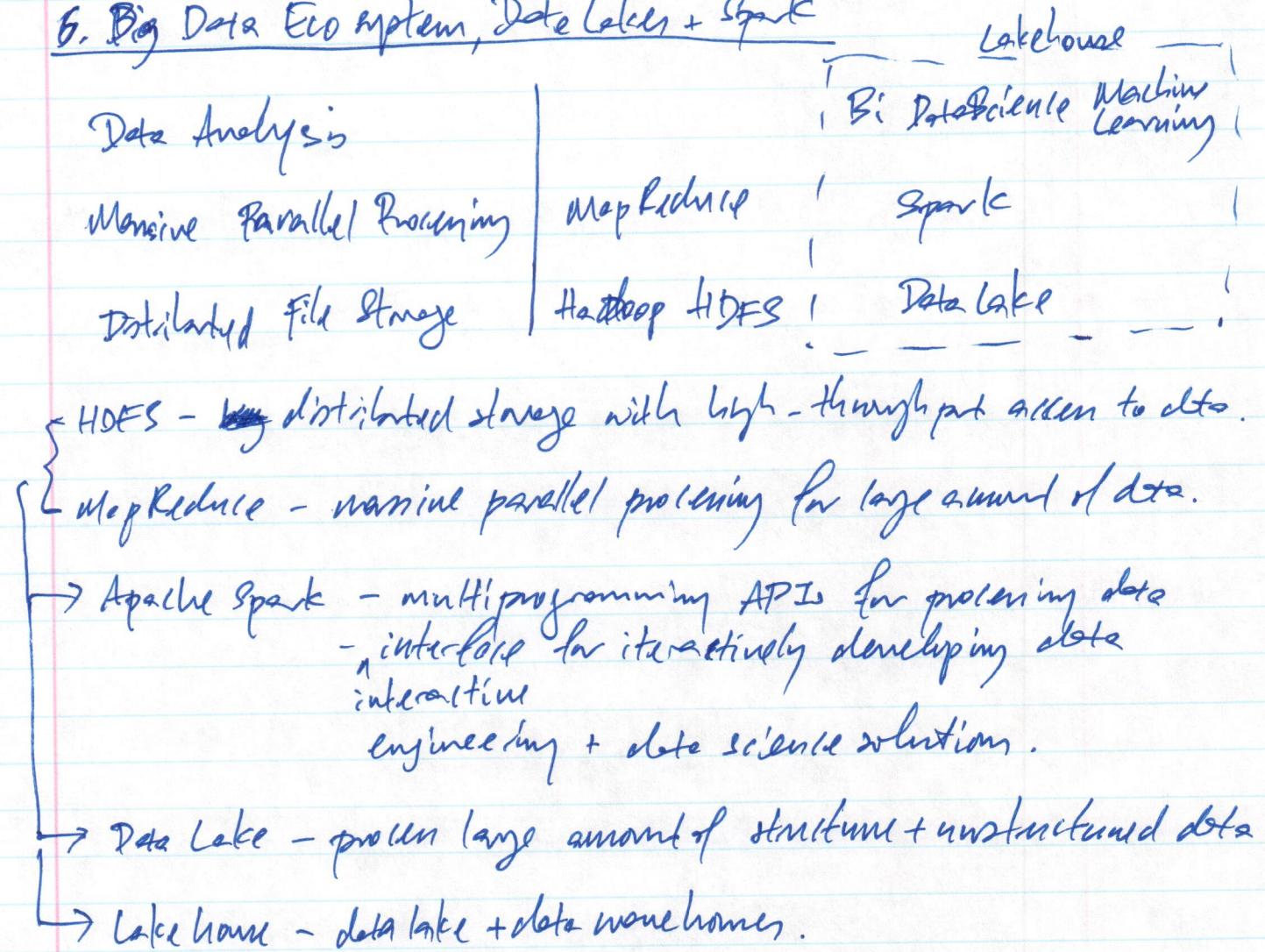
Create Table:

Develop → + Add New Resource → SQL Script → Connect to dedicated SQL pool

CREATE TABLE newTables (
 id VARCHAR(50)
 date DATE
)
Lotto.

INSERT INTO newTables (id, date)
SELECT id, date
FROM staging-table
GO;

6. Big Data Ecosystem, Data Lakes + Spark



From Hadoop to Data Lakehouse:

Bi Reports

Bi Reports DS ML

Bi Report DS ML

Apache
Spark

Data Warehouse

ETL

Structured
Data

Data Warehouse

Spark ETL

Data Warehouse

ETL

Structured, semi-structured,
unstructured Data

Data Lake

Metadata +
Governance layer

```

graph TD
    DL[Data Lake] --> MGL[Metadata + Governance layer]
    MGL --> DLH[Data Lakehouse]
    ETL[ETL] --> MGL
    Spark[Spark] --> MGL
  
```

Structured, semi-structured,
unstructured Data

Data Lakehouse

~~Hadoop~~

- ecosystem of tools for big data storage + analysis.
- writes intermediate results to disk whereas Spark keeps data in memory for as long as possible.
- Spark is faster; does not include file storage system.

Hadoop MapReduce - processing + analyzing large datasets in parallel.

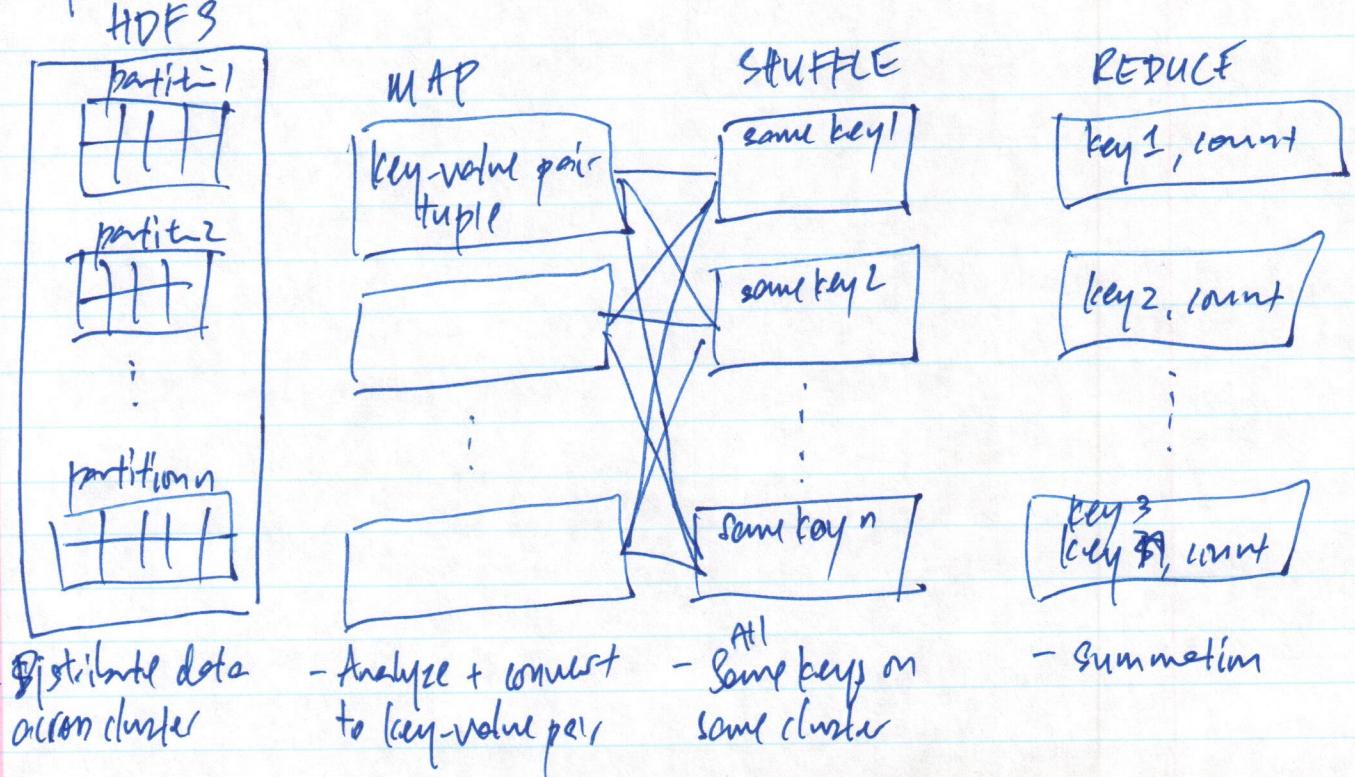
Hadoop YARN - resource manager to schedule jobs across cluster.
what resources are available, assigns tasks.

HDFS - big data file storage system that splits data into chunks + stores chunks across cluster of computers.

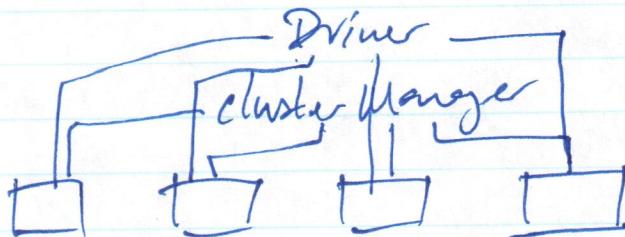
Apache Pig - SQL-like language that runs on top of Hadoop MapReduce.

Apache Hive - another like Apache Pig

MapReduce:

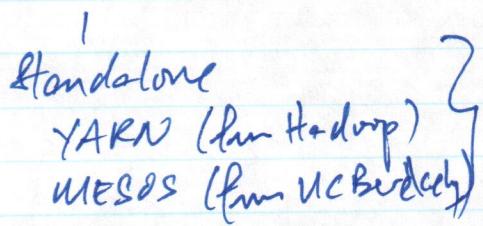


Spark:



→ monitor available resources + make sure all machines are responsive

cluster modes:



3 types of cluster managers

use cases for Spark:

- data analytics
- machine learning
- streaming
- graph analytics

Limitations Spark:

- Spark streaming has latency of at least 500 ms since it operates on micro batches of data, rather than 1 record at a time.
- Other streaming tools - Storm, Apex, Flink can reduce this latency.
- Can't do ML with Spark as it only supports algorithms that scales linearly with data.

Key feature Data Lake: and semi-structured.

- includes unstructured data in addition to structured.
- lower costs with ETL/ELT operations using big data tools.
- schema-on-read over schema-on-write, so lowers costs + work on ingesting large amounts of data.

Apache Spark works on

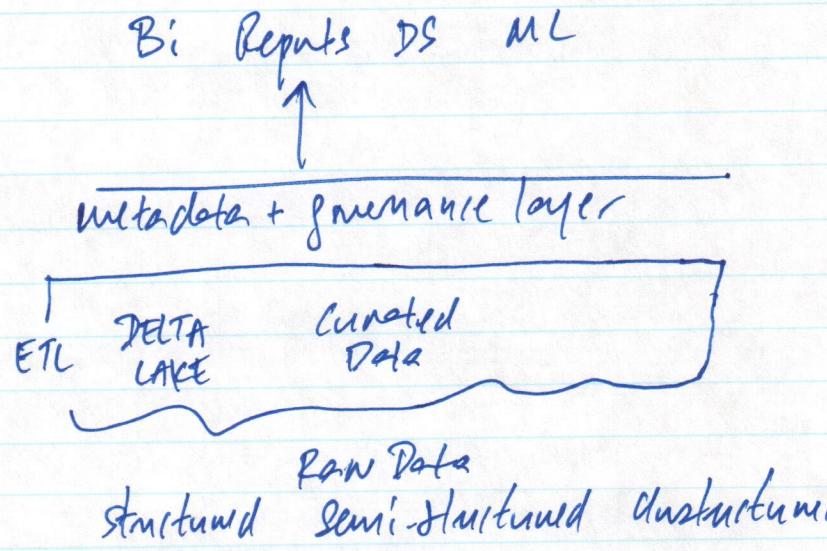
- DE
- Data Analytics
- DS

Data Lake architecture limitation:

- lost atomic transactions
 - ↳ failed jobs left in corrupt state.
- lost quality enforcement
 - ↳ inconsistent + therefore unusable data.
- lost consistency in data structures
 - ↳ impossible to stream + batch process while ingesting.

Data Lakehouse

- addition of metadata + governance layer on top of data lake.
 - ↳ creates pool of raw data + curated data!



- ingest large amounts of data quickly
incrementally improve quality of data.



2. Data Wrangling with Spark

Functional programming - Spark
vs

Procedural programming - C++ / Python

Pure function - functions preserve their inputs + avoid side effects.

Spark - lazy evaluate - build step by step direction

↳ Directed Acyclic Graph (DAG)

- wait till necessary, then run DAGs
- multi-step combos are called stages.

```
import findspark
findspark.init('spark-...-bin-hadoop')
```

```
import pyspark
sc = pyspark.SparkContext(appName = 'example')
```

```
distributed_list = sc.parallelize([1,2])
```

```
distributed_list.map(lambda x: x.lower()).collect()
```

```
distributed_list.collect()
```

↳ not mutated

→ gather result from all machines.

Hadoop + DFS

- splits files into 64 or 128 mb blocks. for digestible chunks.

- stored in fault tolerant way.

Spark Dataframe:

```

import pySpark
from pySpark import SparkContext, SparkConf
conf = SparkConf().setAppName('name').setMaster('ip address')
or
'local'
sc = SparkContext(conf=conf)
from pySpark.sql import SparkSession
spark = SparkSession.builder \
    .appName('example') \
    .getOrCreate()

```

spark.sparkContext.getConf().getAll()

```

log = spark.read.json(pathfile)
log.printSchema()
.log.describe()
.log.show(n=1) ←
.log.write.csv(out-path, 'csv', header=True)

```

spark_df to pandas_df:

pandas_df = spark_df.toPandas()

Typical spark functions:

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

Aggregate functions:

- count()
- countDistinct()
- avg()
- max(), min()
- agg() after groupBy()

User Defined Function (UDF)

Window functions

- partitionBy
- rangeBetween, rowsBetween.

Spark SQL:

```
log.createOrReplaceTempView('table')
```

```
spark.sql("SELECT * FROM table LIMIT 5").show()
```

```
sql_df = spark.sql("""  
    SELECT *  
    FROM table  
    LIMIT 5  
""").collect()
```

```
pandas_df = sql_df.toPandas()
```

UDF:

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

```
spark.sql("""  
    SELECT *, get_x(x) AS x  
    FROM table  
    LIMIT 1  
""").collect()
```

Resilient Distributed Data Set (RDD)

└ <1.3 RDD - low level abstraction of data

spark version - 1.3 DataFrame API

└ 2.0 Data Frame Data Sets API

Spark Accumulators

└ 7. Spark Debugging and Optimization

x = SparkContext.accumulator(0, 0)

Spark Broadcast

Action Functions

- trigger lazy evaluation functions
- e.g. .write()

Data Screening

- check min, max, ranges, value-count()
- change workload division - create composite keys
- break data into smaller partitions.
- ↳ df.repartition(number_of_workers)

q. Azure Data Bricks

Azure Databricks contains:

- distributed filesystem
(storage)
- workspace (data engineering, data analytics, data science)
(ingesting, processing, transforming, storing, serving data)
- spark clusters
↳ Spark API, Spark SQL, Spark DataFrames

User of Data Bricks:

1.

Unstructured /
Structured Data
→ Azure
Data
Factory
(ingestion)

→ Azure
Data Lake
Storage

→ DataBricks
(Spark
processing)

↗ Azure
Synapse
Analytics

↗ Power BI

2.

Azure SQL Database

IoT Hub

Stream Analytics

(data flow
processing)

DateBricks
↓
ML

→ Data Lake → Power BI

↗ DataBricks
↓
Data Lake

Azure Portal → Create Resource → Azure Data Bricks

↳ Name
Resource Group
Region
Pricing - Standard

Launch Workspace

↳ Compute → Create Cluster

↳ Name

Cluster Mode
Databricks Runtime
Autopilot Options

Worktype

↳ Min Workers, Max Workers → determine cost

Driver Type

→ Create Cluster

→ Stop / Start / Delete

→ ~~Workspace~~ → Shared

↳ Create

↳ Notebook

↳ Name

language
cluster

cluster must run first.

Write To Delta Lake:

df = spark.read.format('csv').option('header', 'true')

display(df)

df.write.format('delta').mode('overwrite').save('filepath')

df.write.format('delta').mode('overwriten')." " . saveAsTable('dataTable')

↳ Create Table

Data → Database Table → default → Table
↳ data table

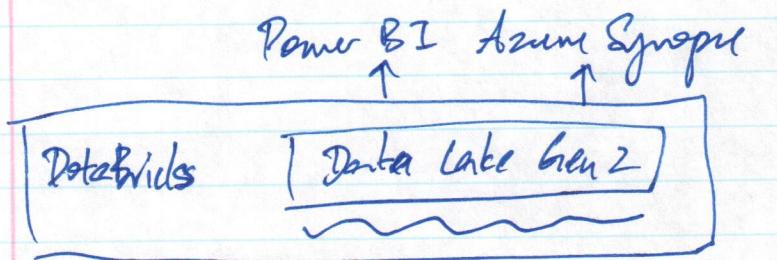
→ DBFS → /delta
↳ delta → -delta-lay
part.....parquet

Settings workspace settings
↳ Admin Console ↳ Advanced

↳ DBFS File Browser - enabled.
(Refresh browser)

10. Data Lakes and lakehouse with Azure DataBricks

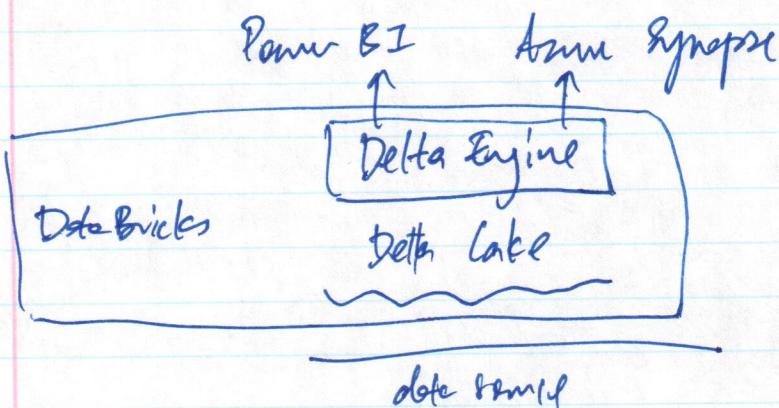
Data Lake Gen 2:



structured, semi-structured
unstructured data

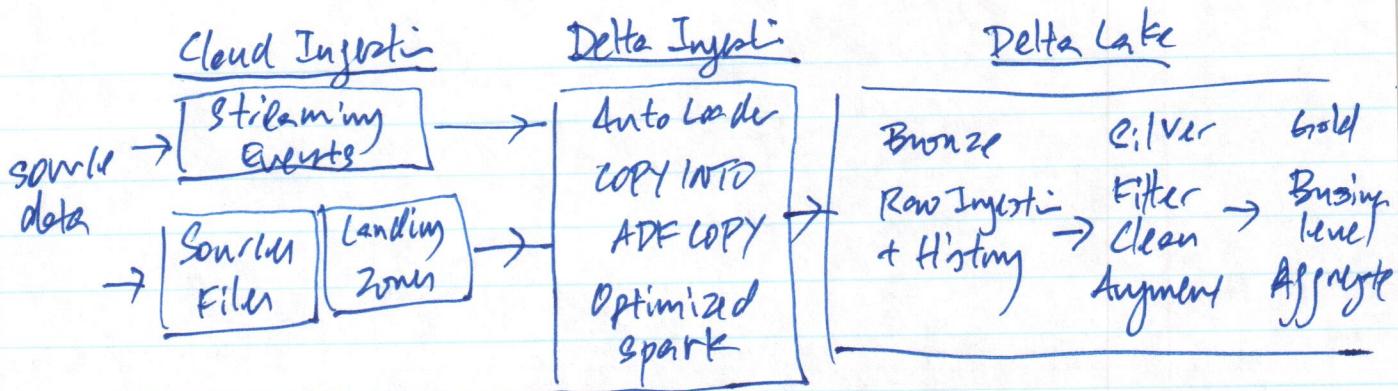
- adds Databricks Storage
- hierarchical namespaces to enable better organization of information.
- actionable using Flink API's
- integrates with Azure Synapse, Data Factory, DataBricks.

Delta Lake:



tasks:

1. create workspace
2. ingest data
3. create, update, delete Table
4. read, write Data



Create Table:

```
spark.sql('CREATE TABLE table (' +
    'loan_id BIGINT,' +
    'paid_amnt DOUBLE')')
```

```
or) df.write.format('delta') \
    .mode('avrowrite') \
    .saveAsTable('table')
```

```
, spark.sql("CREATE TABLE datatable USING DELTA \
    LOCATION '/delta/data'")
```

Drop Table:

```
spark.sql("DROP TABLE IF EXISTS datatable")
```

Read data:

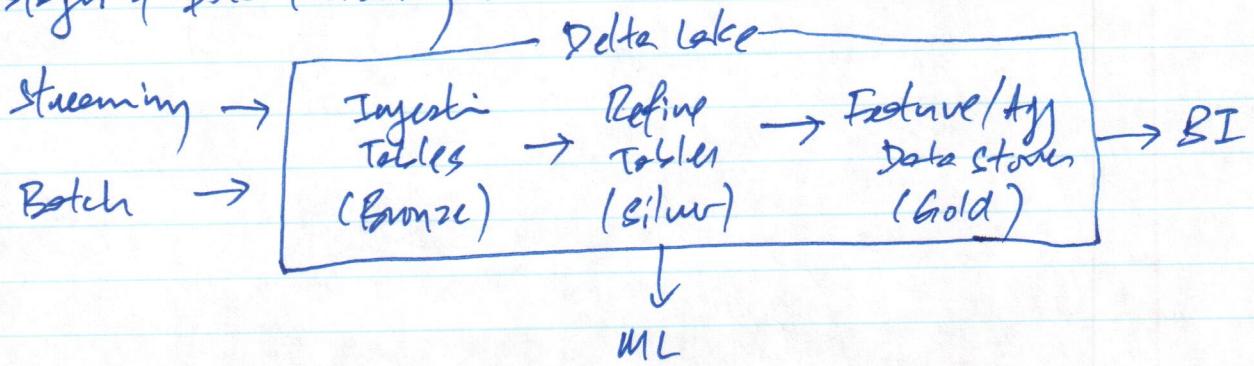
```
spark.table('default.tablename')
spark.read.format('default').load('/delta/tablename')
```

Write data:

```
df.write.format('delta').mode('append') \
    .saveAsTable('default.tablename') \
    .save('/delta/tablename')
```

```
spark.sql("CREATE TABLE tableone")
    "USING DELTA LOCATION '/delta/data'")
```

Stages of Data Processing:



Ingestion:

```
df = spark.read.format('csv') \
    .option('inferSchema', 'false') \
    .option('header', 'true') \
    .option('sep', ',') \
    .load('Filestone/demo2/file.csv')
```

```
df.write.format('delta') \
    .save('delta/bronze-file')
```

Refine:

```
df = spark.read.format('delta') \
    .load('delta/bronze-file')
```

from pyspark.sql.functions import sha2, concat_ws

df = df.dropDuplicates(df.columns)

df = df.withColumn('hash1', sha2(concat_ws('||', *df.columns), 256))

```
df.write.format('delta') \
    .mode('overwrite')
```

```
.saveAsTable('silver-table')
```

Aggregate:

```
df = spark.read('silver-file')
```

```
df = df.groupBy('x', 'y').count()
```

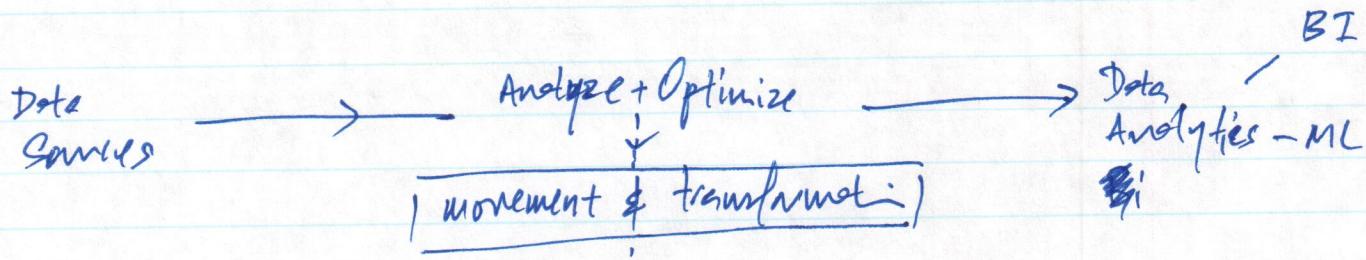
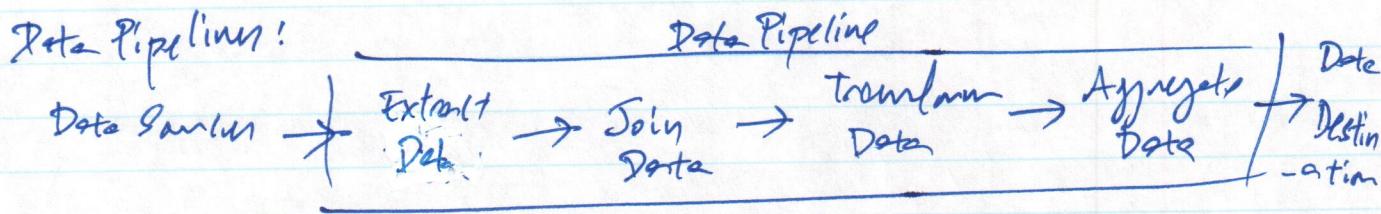
```
df.write.format('delta')
```

```
.mode('overwrite')
```

```
.saveAsTable('gold-file')
```

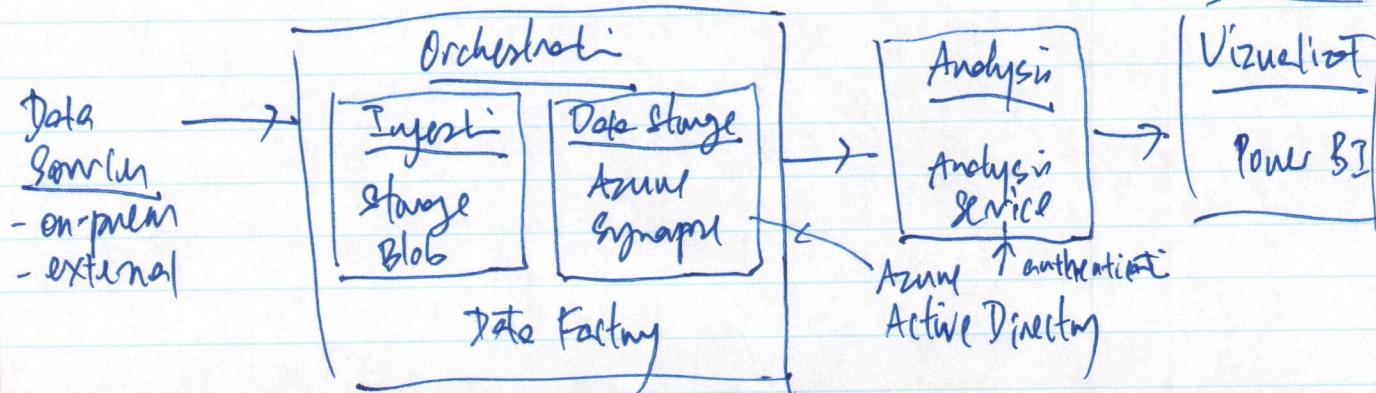
11. Data Pipelines with Azure

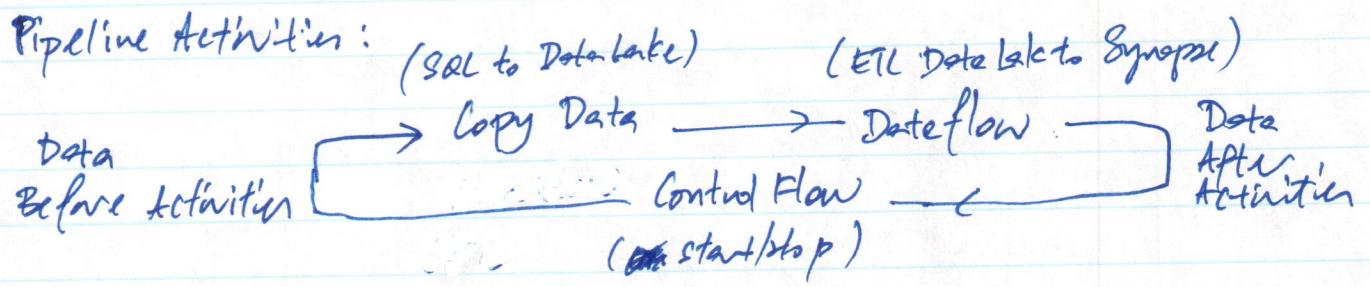
Data Pipelines:



Azure Data Factory \longleftrightarrow Azure Synapse

- provide cloud-based code-free ETL or ELT as a service to orchestrate data movement





Azure Portal → Create Azure SQL Database

- ↳ Query Editor
 - ↳ login
 - ↳ query

→ Create Azure Synapse

↳ Open Synapse Studio

↳ Manage - SQL pools - New

↳ (Create)

↳ Develop - New SQL script

↳ Create tables in SQL pools

→ Create Azure Data Factory

↳ Name
Region
Version

↳ Git (configure later)

↳ Create

↳ Open Azure Data Factory Studio

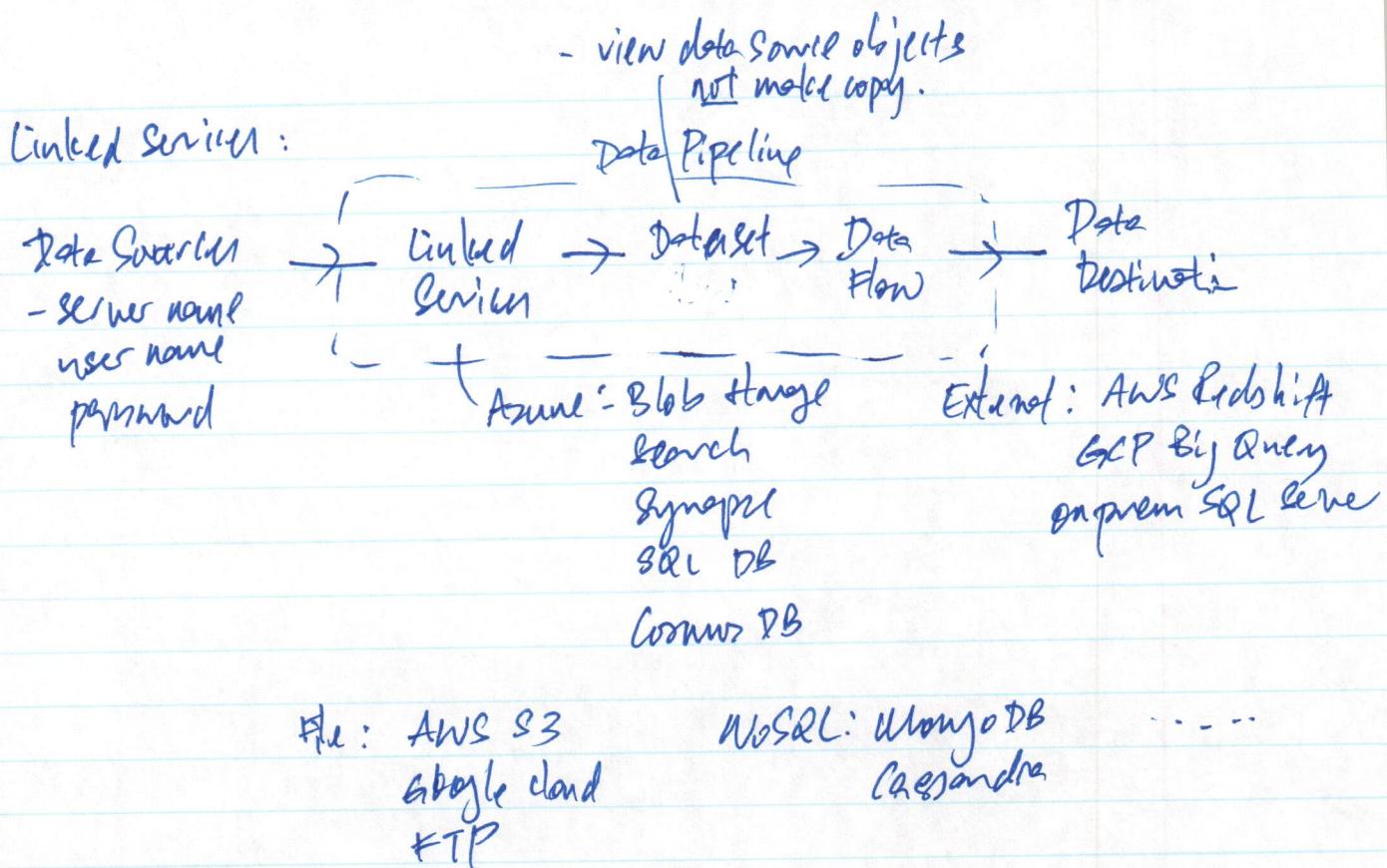
→ Create Storage Account

↳ Advanced

↳ Enable hierarchical namespace

↳ (next)

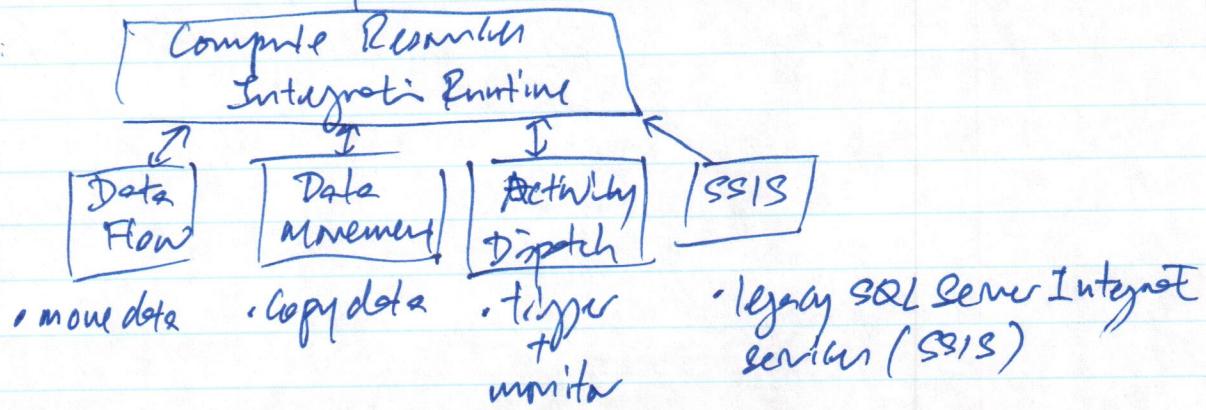
↳ Storage Browser → Blob container - + add container



↳ Azure Data Factory

- ↳ Manage - Linked services
 - SQL (white list IP address)
 - Synapse
 - Storage account
- ↳ Author - Data Sets - New Data Sets
 - ↳ SQL
 - ↳ Name
 - ↳ linked service
 - ↳ table name
 - ↳ Azure Synapse Analytics
 - ↳ name
 - ↳ linked service
 - ↳ table name
- ↳ Publish

Integrate Runtime (IR) :



3 types:

1. Azure IR

- ↳ data flow
- ↳ data movement } between cloud stores
- ↳ dispatch to external compute
- ↳ SSIS in public or private network
- ↳ serverless compute.

2. Self-hosted IR

- ↳ data movement between cloud stores and private network.
- ↳ dispatch to external compute ~ Azure Virtual Network.

3. Azure SSIS IR

- ↳ lift + shift existing SSIS to execute in Azure.

Azure Data Factory

↳ manage → Integrate Runtimes - + New

↳ Azure, self-hosted

↳ Azure

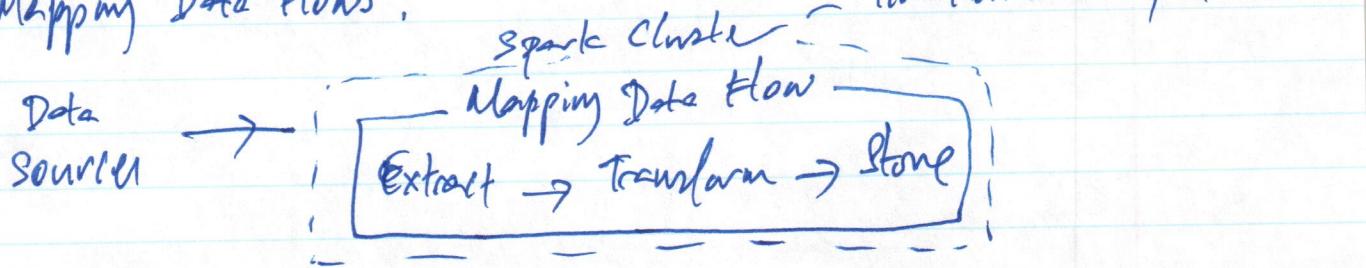
↳ cloud, region

↳ Time to Line

Compute type, cloud

12. Transforming Data in Azure Data Pipelines

Mapping Data Flows:



3 types:

1. Schema modifiers

- create new derived columns based on calculations, aggregates data,
or pivoting data

2. Row modifiers

- change rows e.g. filtering, sort, filter, based on insert/update/delete/upsert policies.

3. Multiple inputs/outputs

- generate new data with joins, unions, or splitting data.

- flowlets - prevent logic duplication + apply consist logic
- formatters - change one format to another

Data flows

- interpreted with Visual Expression Builder in Azure Data Factory
to perform logic transformation as simple expression.

- Expression - column mapping } evaluated Spark data types
- functions
- parameters } at runtime.

+ - input schema - reference columns from input data source.
functions - aggregate, convert, date + time -
parameters - parameters passed.
cached lookup - lookup on data.

Arrange Date Factor -

- Author - Data flows - New Data Flow
 - Add Source
 - Name
 - Select - Dataset
 - Data flow delay
 - Spin up Spark cluster
 - Delay time to line
 - Debug setting
 - Concurrent
 - Data Preview
 - Refresh

> Destination

- sink
 - Name
 - Dataset in Synapse
 - Mapping
 - Auto mapping - uncheck to confirm
 - Data Preview

Arrange Synapse Analytics

- Develop - New SQL Script

- CREATE TABLE table

ADF

- Author - Datasets - New Dataset

- Synapse

- Name

- linked XML

- tablename

- Data flow - New

- Add Source

> filter - name

filter on - Open expression builder
Data Preview - refresh
Save + Finish

> Derived - name

Column column name expression builder
Save + Finish

> Sink - name

dataset - synapse table

Publish data flow.

'Author - Pipeline - new'

'Drop Data flow'

'Name'

'Setting - Data flow name'

- Step my linked service
Step my storage folder

- Validate

- Publish

- Add trigger

trigger now or schedule time

'Monitor - pipeline runs refresh'

Synapse → confirm data in table

~~#FF~~

~~Author - Pipeline - new~~

Monitor A pipeline

- get alerts

- Run monitor

ADF - pipeline - Delay

- ` with data flow runtime
- ` use delay cluster

Details - stages - lineage

Anne Synapse trinity

- ` manage - linked service

- ` data - linked

- ` synapse storage account

- ` upload - csv

- ` right click - load to dataframe

- ` manage - spark pool - new

- ` name

- ` max size, number

- ` develop - ~~notebook~~

- ` attach to ~~spark~~ pool
- ` read csv
- ` create database table
- ` save as table

- ` integrate - new pipeline - drag notebook

name

setting - select notebook

spark pool

- ` validate
- ` publish
- ` trigger

- ` monitor
- ` refresh

Power Query: (code-free, excel like transformations)

- uses M language
- ADF uses Apache Spark behind scenes, automatically converts M language types to Spark.

ADF

 \ Author - Power Query - New

 \ Select Data set
 (Action similar to Excel)

 \ ---
 - Publish

 - Pipeline - new - Power Query (adrey)

 \ new

 \ settings - select , staging

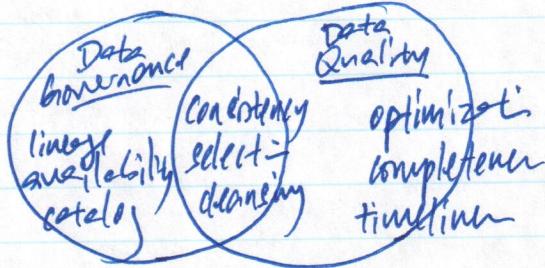
 \ Sink - select - db.

 - validate

 - publish

 - trigger

13. Azure Pipeline Data Quality



Data Governance

- track data lineage

- know who it's available to

- understand catalog of data available

Data Quality

- optimizing data so it's complete, timely, consistent.

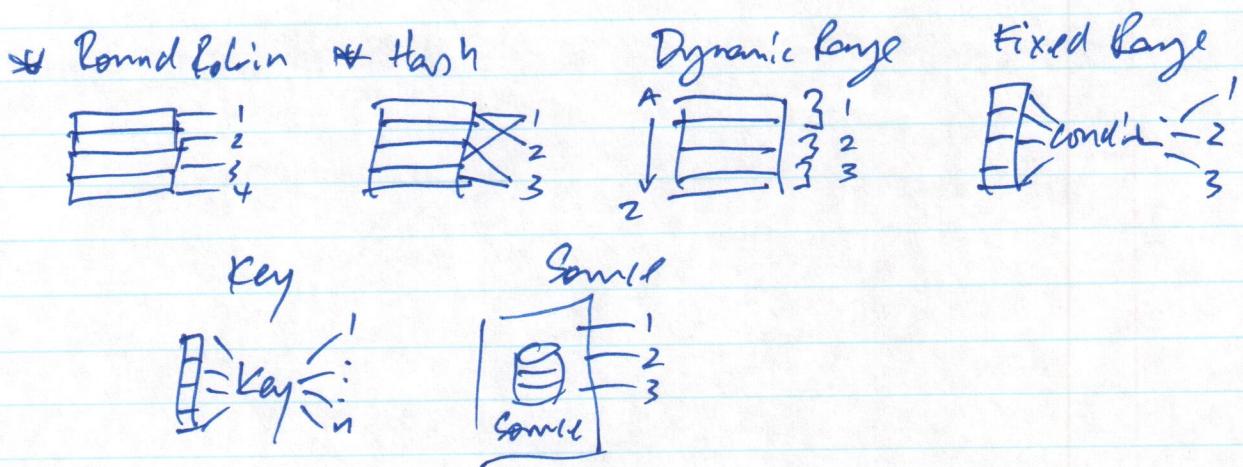
Best practices to ensure ADF or Synapse achieve minimum scale

1. Tuning of Integratiⁿ Runtime

- ↳ Spinning cluster for Spark compute takes 3-4 minutes if concurrent jobs, make sure 'Time to live' is long enough, subsequent

2. Partitioning data

1. Current partition - recommended for all except SQL db.
2. Single partitⁿ - smaller dataset
3. Set partition count - recommended for SQL db.



3. Schema Drift

- ↳ Columns in source data changes after Pipeline is developed + in use.
- ↳ enable "Allow schema drift" in Data flows for source + sink data xt.

Slowly changing dimensions

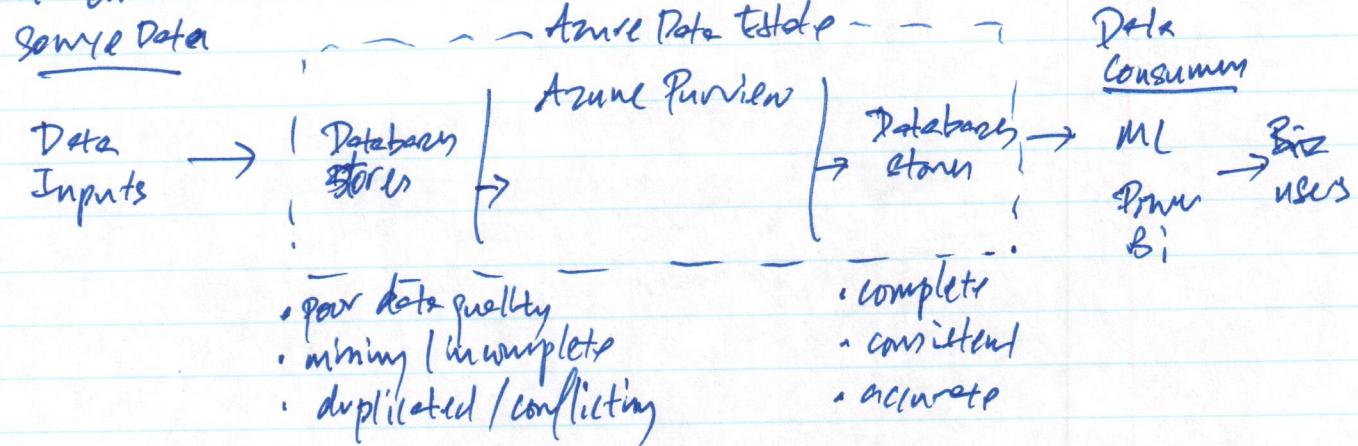
- ↳ e.g. customer phone, address, product name.
- ↳ 6 types of SCD
 - 0: ignore changes, keep only original
 - 1: overwrite existing with new
 - 2: add row with new values, add version column
 - 3: add new columns for new values
 - 4: maintain all historical values in new history table.

5 : extension of Typl, where a multi-dimension table is created by maintaining keys.

6 : combination of 1, 2, 3 - adds new record for new values + maintains a new column,

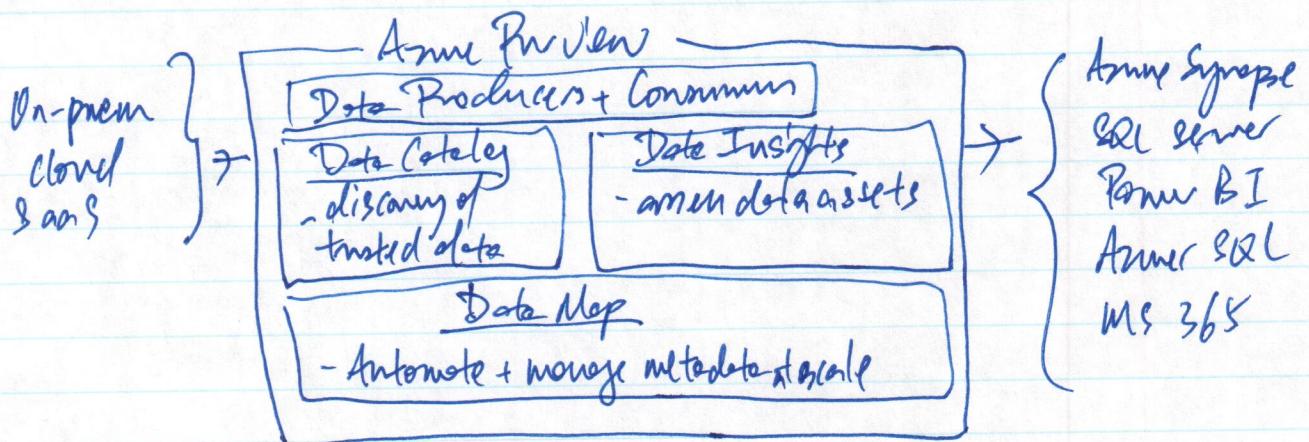
7 :

Data Governance:



Azure Purview -

- data governance service
- discover data quickly, classify data, data lineage
- creates data map of data estate with datasets + their relationships + stores the corresponding metadata.
- connects with ADF, Synapse.



Create Purview account

ADF

↳ manage - Azure Purview - Connect to Purview account
↳ authn - pipeline - trigger now

Azure Purview - Open A3 Purview studio

↳ Browse assets
↳ collection
↳ SQL
↳ Lineage

4. Azure Data Pipelines in Product:

Parameterizing Pipelines:

- supported by pipeline, datasets, linked services, dataflows.
- examples:
 1. linked services - database name, username
 2. pipeline - parameters passed from pipeline to dataflow.
 3. global - used by all pipeline, a dataflow
 - Dev to QA to Prod env. pass to
 - only for ADF, not Synapse
 4. System variables - `(@pipeline().DataFactory())`
`(@pipeline().RunId)`
`(@trigger().startTime, ...)`

ADF

↳ DataFlow - click on canvas
↳ Parameters - New
↳ name, type, value
use Parameter in Open Expression Builder

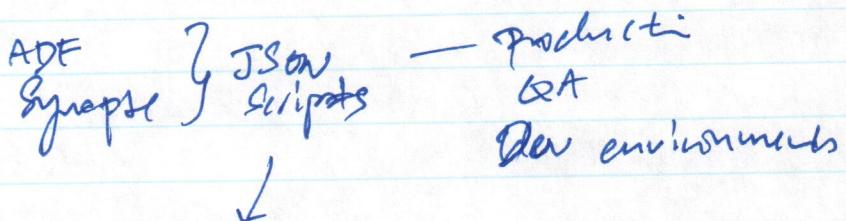
↳ Pipeline - Some on Siffler to create

- to pass to Dataflow, ~~select Dataflow_Parameter~~
↳ click on value - select Pipeline expression

Select Parameters

~~Manage~~ → Global Parameters - New
(available in ADF or Synapse)

- Linked Service - Select - Parameters - New



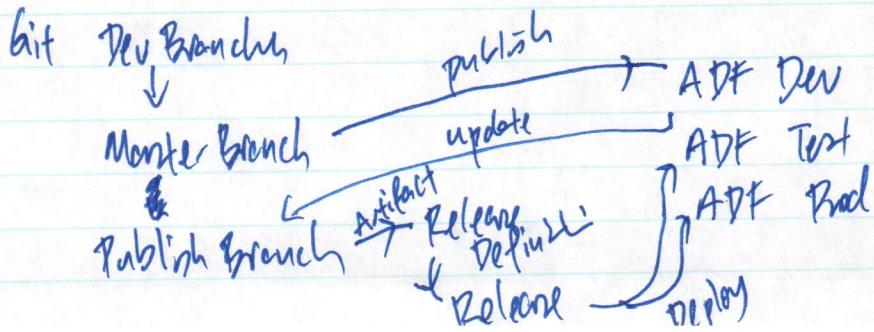
Create with - CLI
PowerShell
.NET
Python
REST
ARM templates

az datafactory create --resource-group ResourceGroupName
--factory-name FactoryName

linked services - {} - copy JSON
Data flow - {} " "

CI/CD:

- integrate ADF, Synapse with Azure DevOps or GitHub
- use ARM templates for target environments



Create project in Azure DevOps Repo.

ADF

↳ manage - Git Configuration - Configuration

↳ Azure DevOps
GitHub

↳ Azure AD

↳ Repo select
project name .
branch name

Author - all Pipelines

Datasets
Dataflows
Link service

Triggers

Integration Runtime

} publish to upo

Project NYC payroll:

