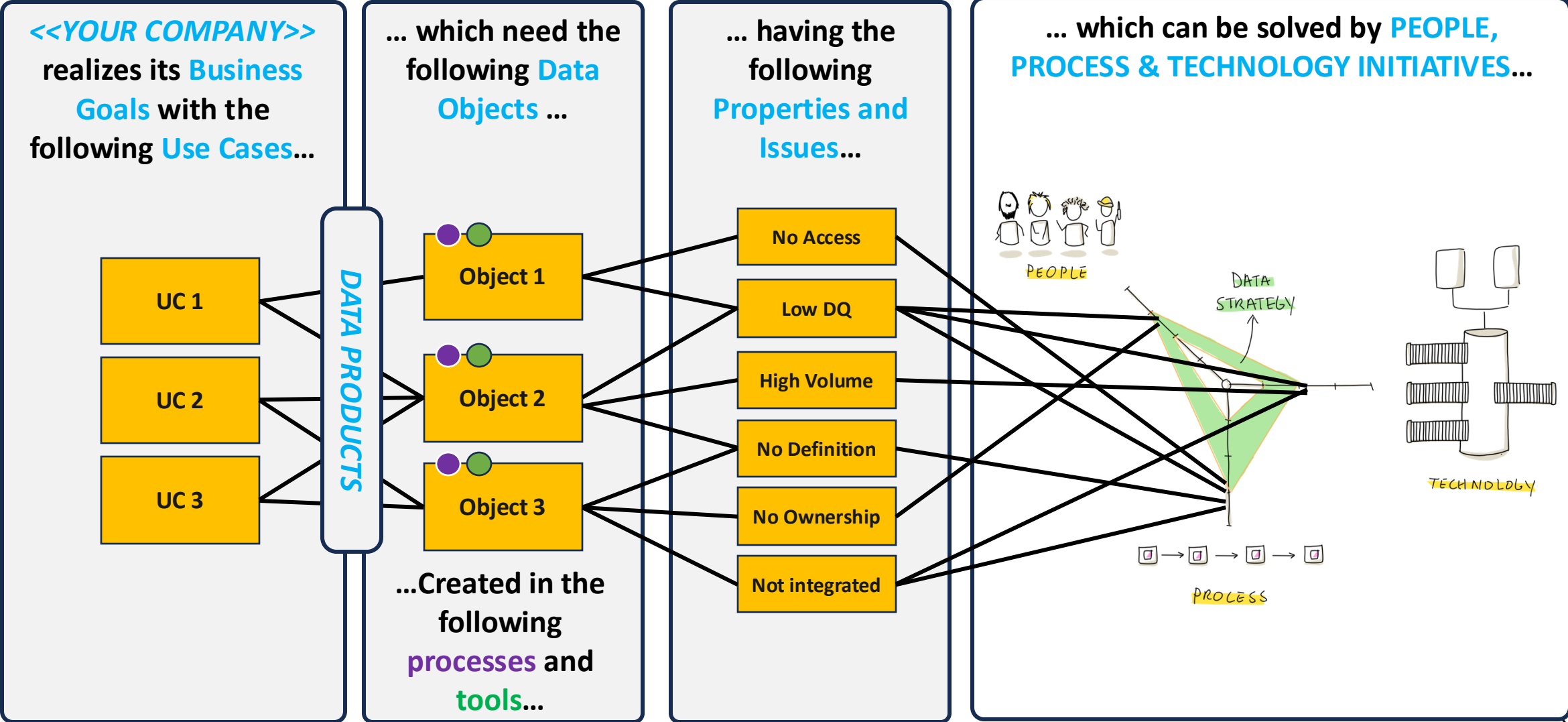


SOLUTIONS

TECHNOLOGY

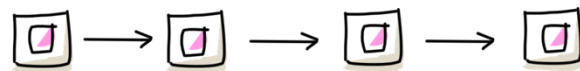
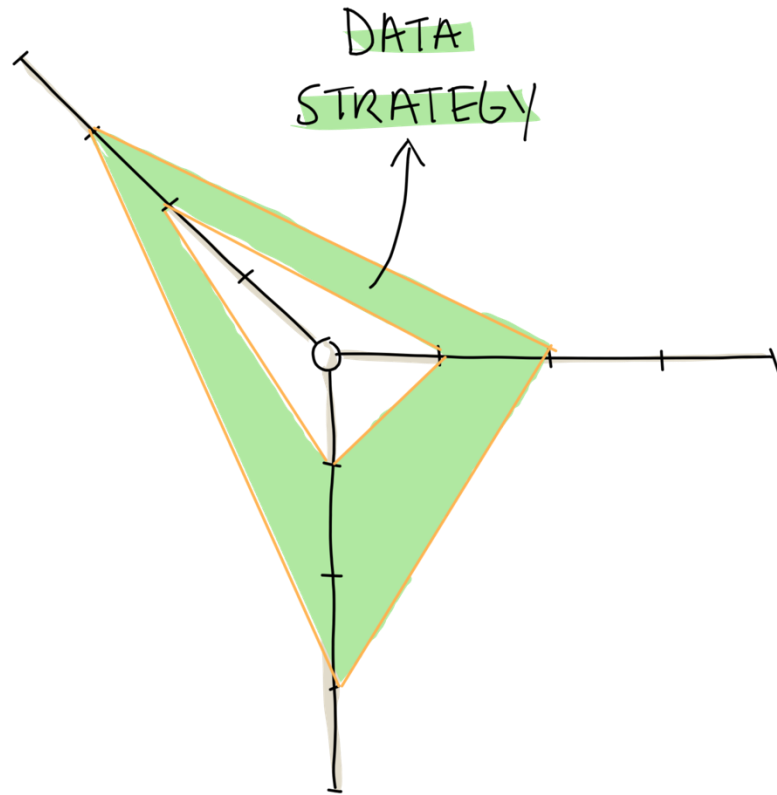


Data Strategy Framework

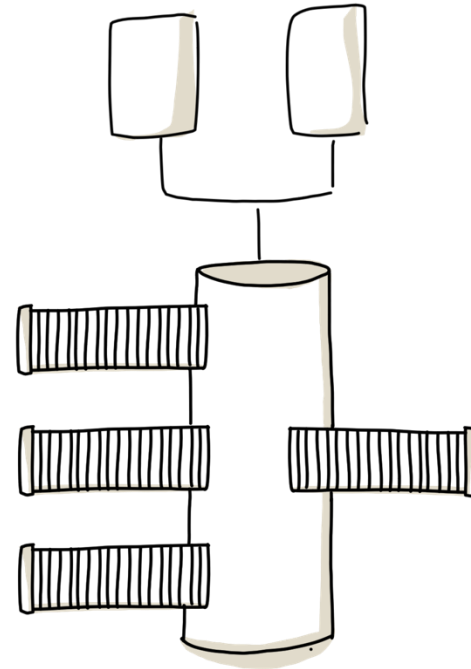




PEOPLE



PROCESS



TECHNOLOGY



Table of Contents

- Dead Horse Theory
- Data Platform
 - Introduction
 - Core Layers
 - Additional Layers
- Technology Selection



Table of Contents

- **Dead Horse Theory**
- Data Platform
 - Introduction
 - Core Layers
 - Additional Layers
- Technology Selection



BEFORE WE DIVE INTO TECHNOLOGY

DATA PLATFORM projects can become LONG
RUNNING & EXPENSIVE projects

BE AWARE OF DEAD HORSES

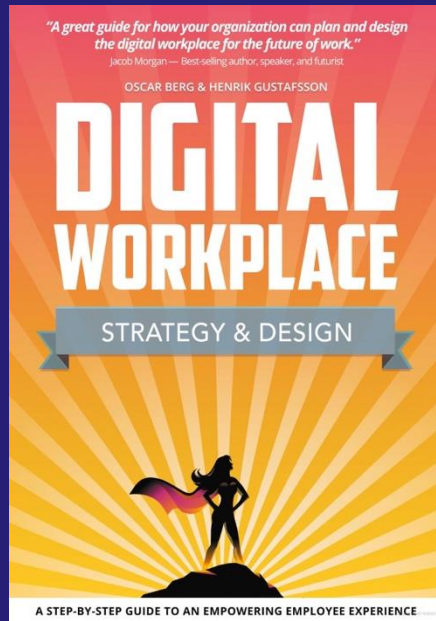
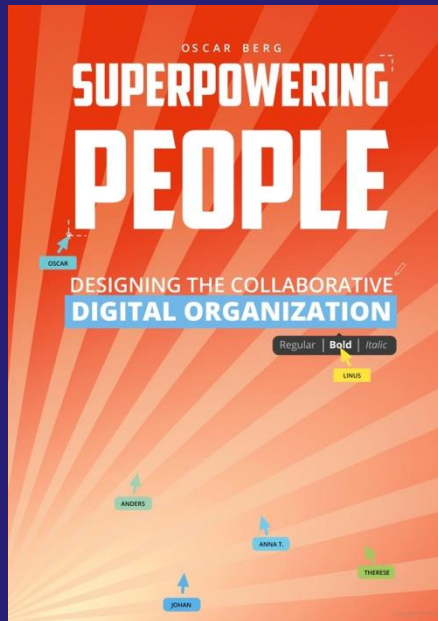
THE DEAD HORSE THEORY *ILLUSTRATED*

This should work!



Source: Oscar Berg

Author of:



Illustrations are part of his new book (to be announced)

The tribal wisdom of the Indians, passed on from generation to generation, says that, "When you discover that you are riding a dead horse, the best strategy is to dismount."

**The Dead Horse Theory
goes on to say
that in modern business,
education and government,
far more advanced
strategies are often
employed, such as:**

1. Buying a stronger whip.

Come on now, move
you stupid horse!



2. Changing riders.

Good luck! Thanks!



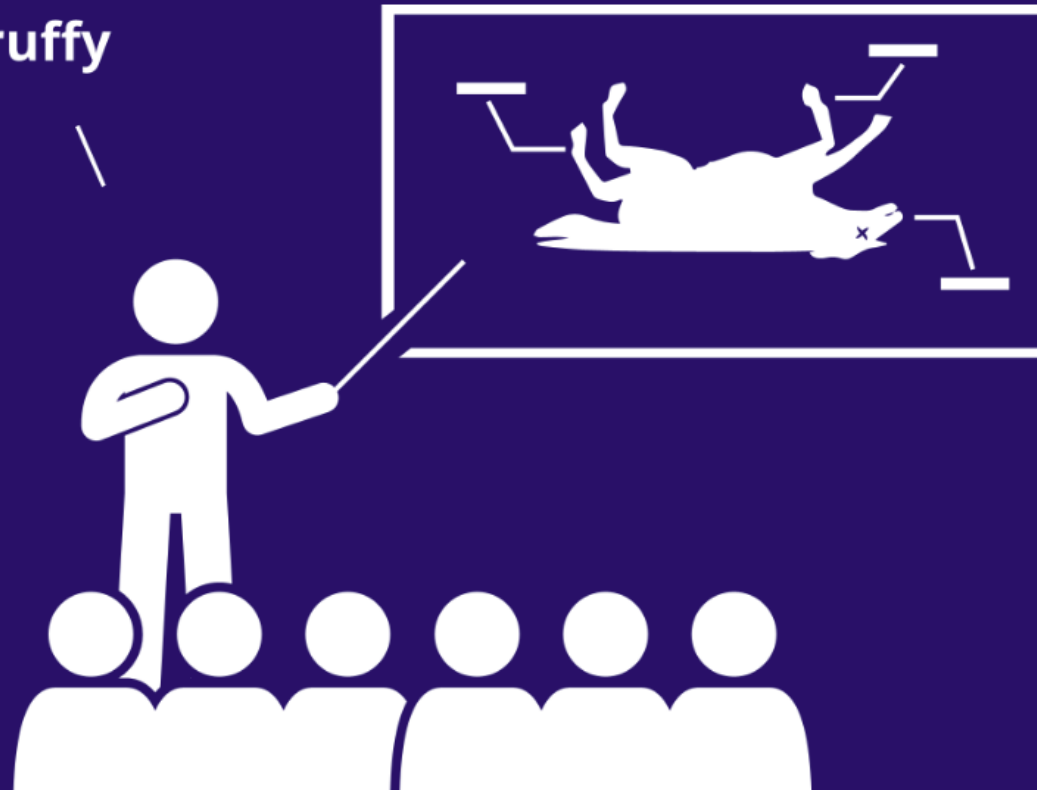
3. Threatening the horse with termination.

There's the door
if you don't
shape up soon!



4. Appointing a committee to study the horse.

The tail seems a
bit scruffy



5. Arranging to visit other countries to see how others ride dead horses.



6. Lowering the standards so that dead horses can be included.

It seems to keep the pace well



7. Re-classifying the dead horse as 'living-impaired'.

This might help



8. Hiring outside contractors to ride the dead horse.

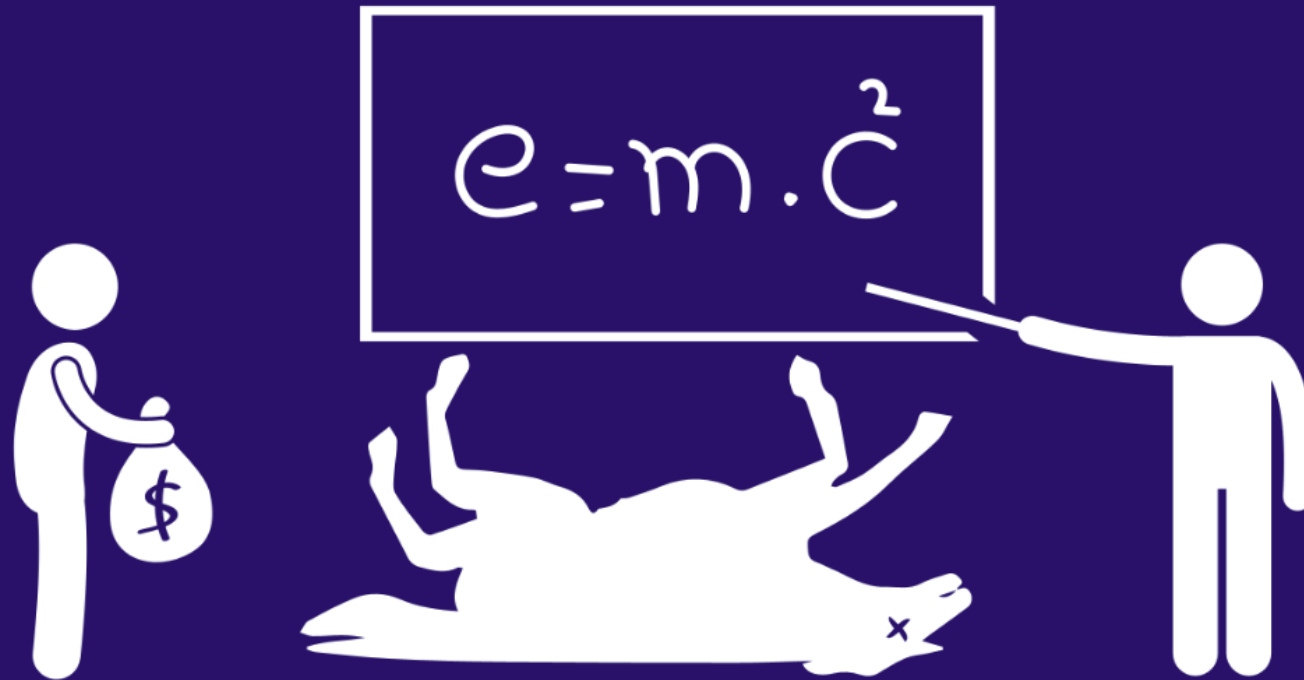
Sign the 1 billion contract here...



9. Harnessing several dead horses together to increase the speed.



10. Providing additional funding and/or training to increase the dead horse's performance.



11. Doing a productivity study to see if lighter riders would improve the dead horse's performance.

Normally we don't hire children, but since you're not being paid it's not considered child labour



12. Declaring that as the dead horse does not have to be fed, it is less costly, carries lower overhead and, therefore, contributes substantially more to the bottom line of the economy than do some other horses.



13. Re-writing the expected performance requirements for all horses.

All we require from you is to be present at the office

Yee-haw!

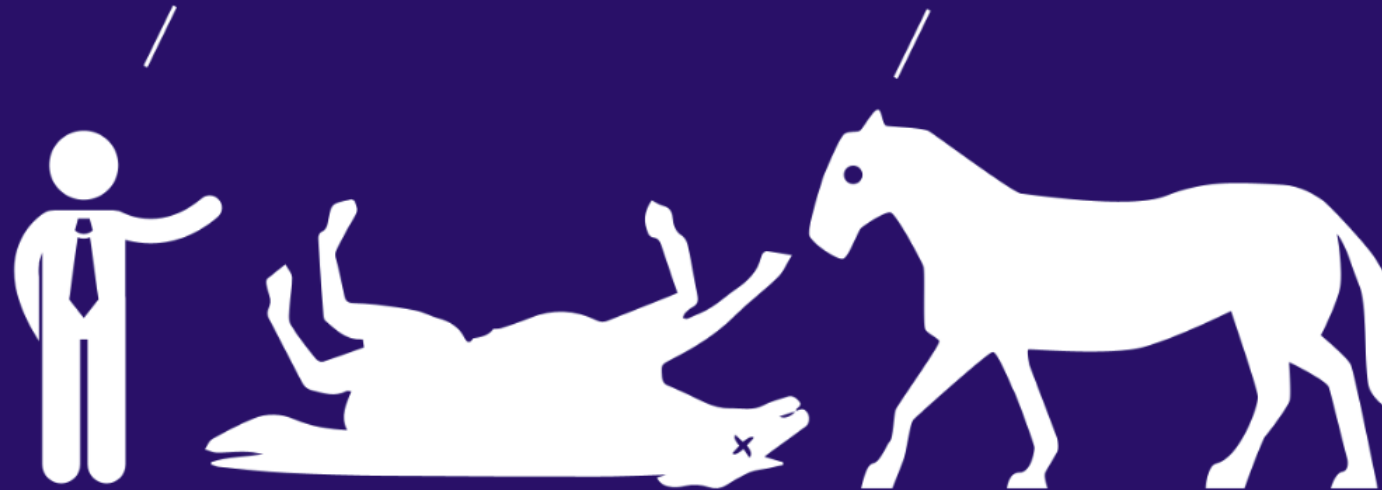
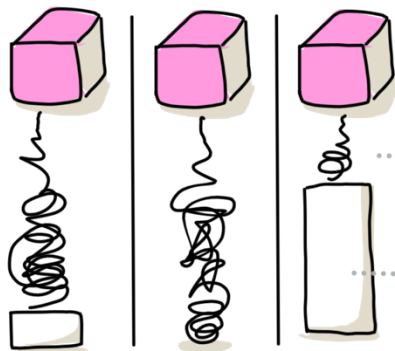


Table of Contents

- Dead Horse Theory
- Data Platform
 - **Introduction**
 - Core Layers
 - Additional Layers
- Technology Selection

MATURITY

HETEROGENEOUS &
UNMANAGED TOOLS

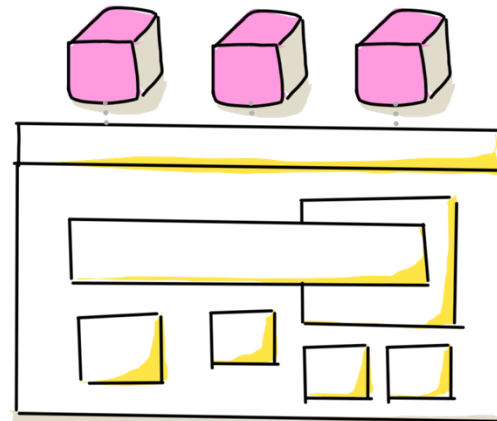


USE CASE

AD HOC

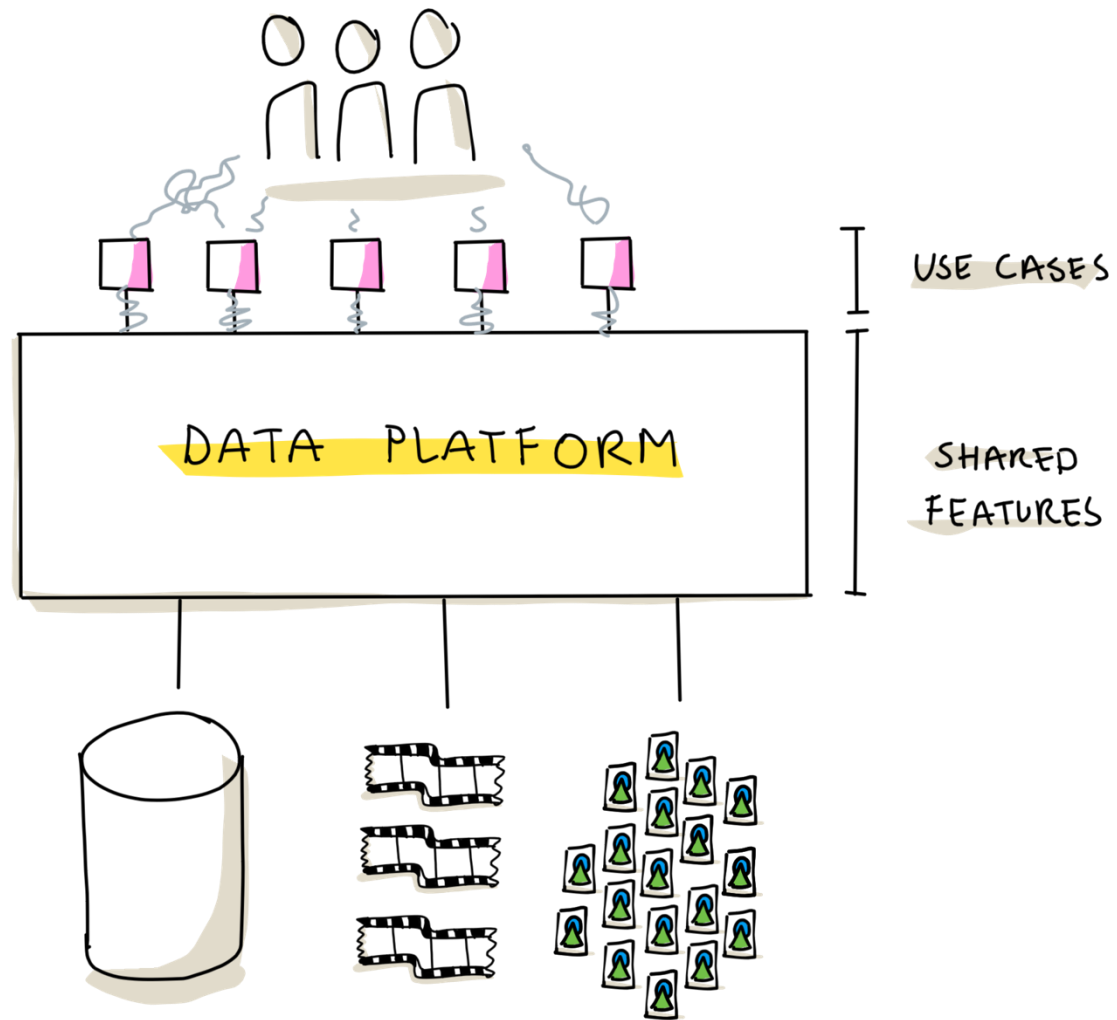
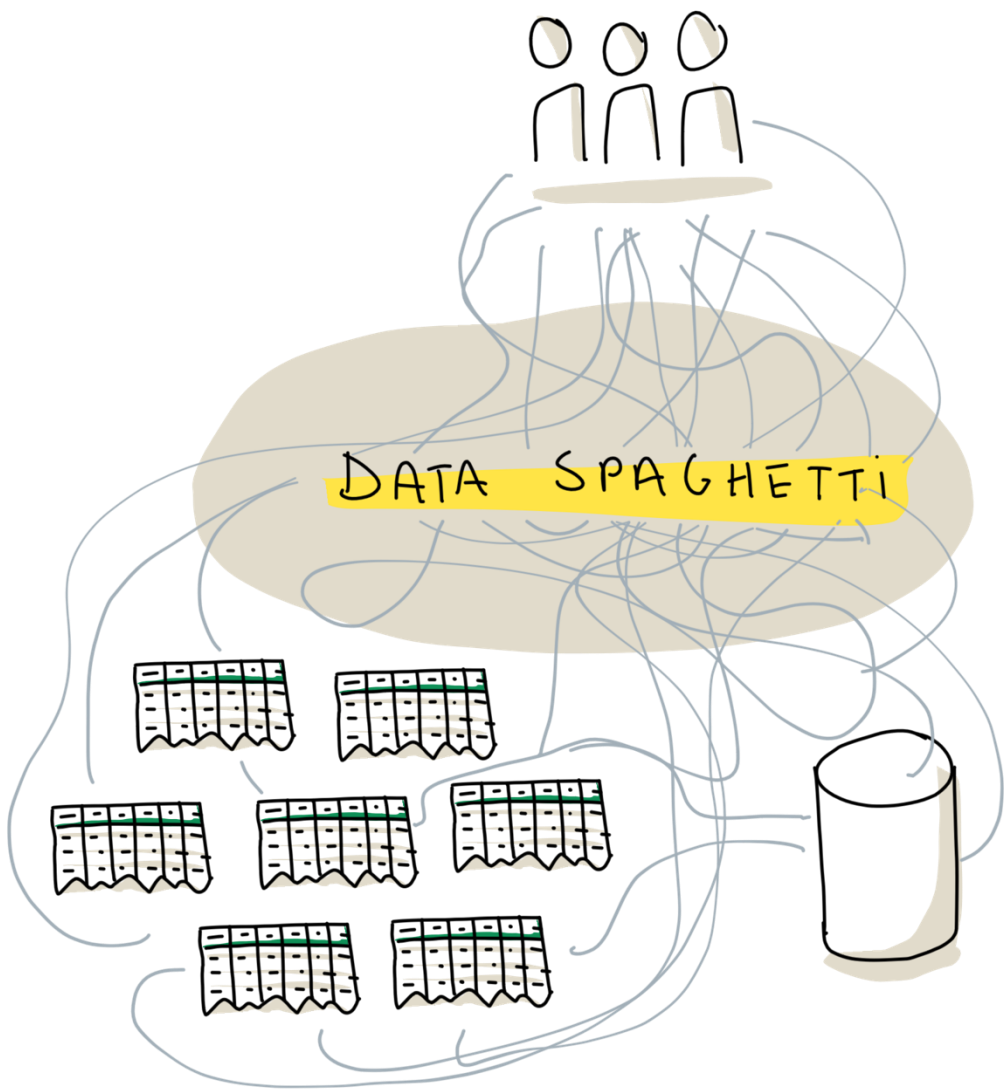
AUTOMATED

DATA TOOL
ECO-SYSTEM

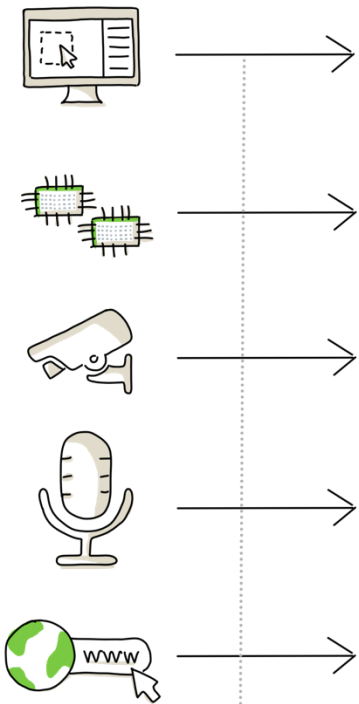


TIME

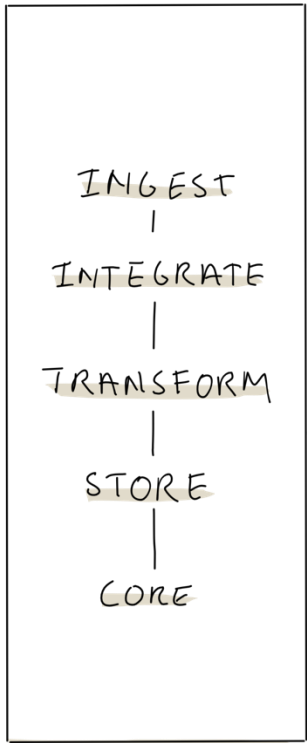




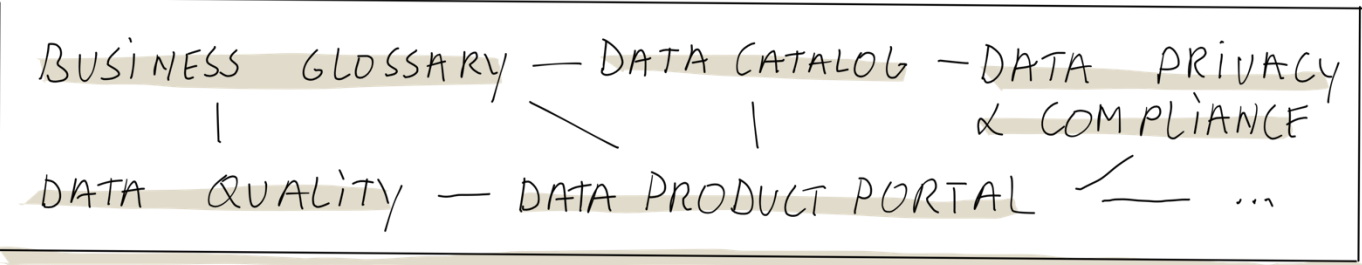
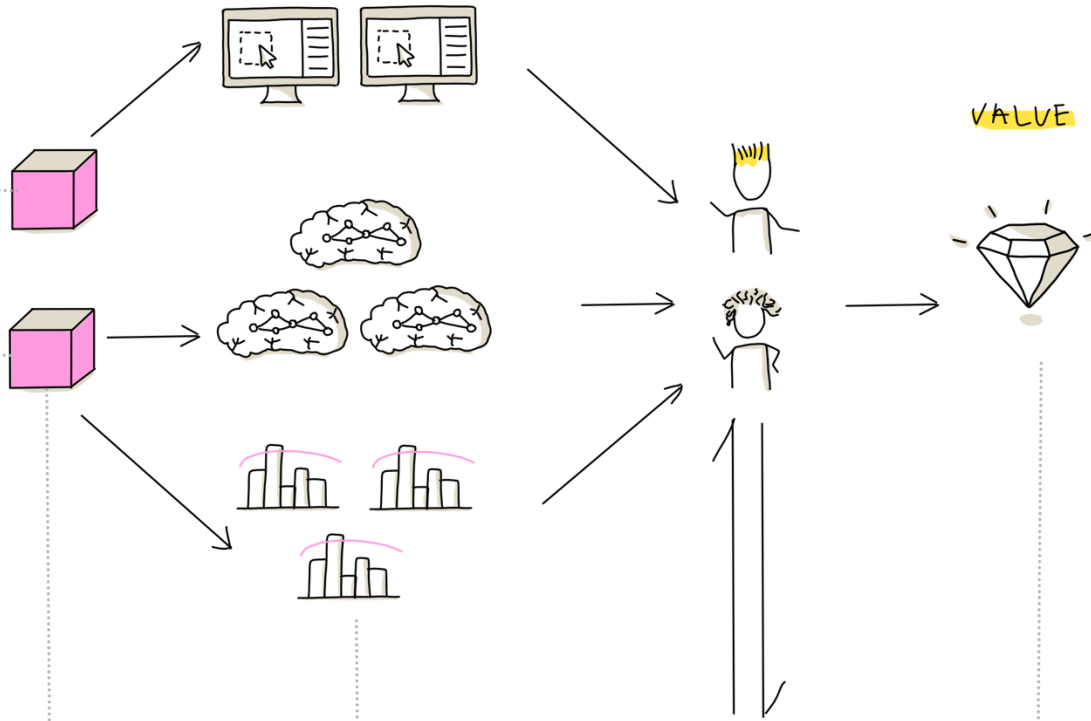
DATA PRODUCERS (SOURCES)



DATA PLATFORM



CONSUMPTION



META-DATA MANAGEMENT



Modern Data Platform Architecture

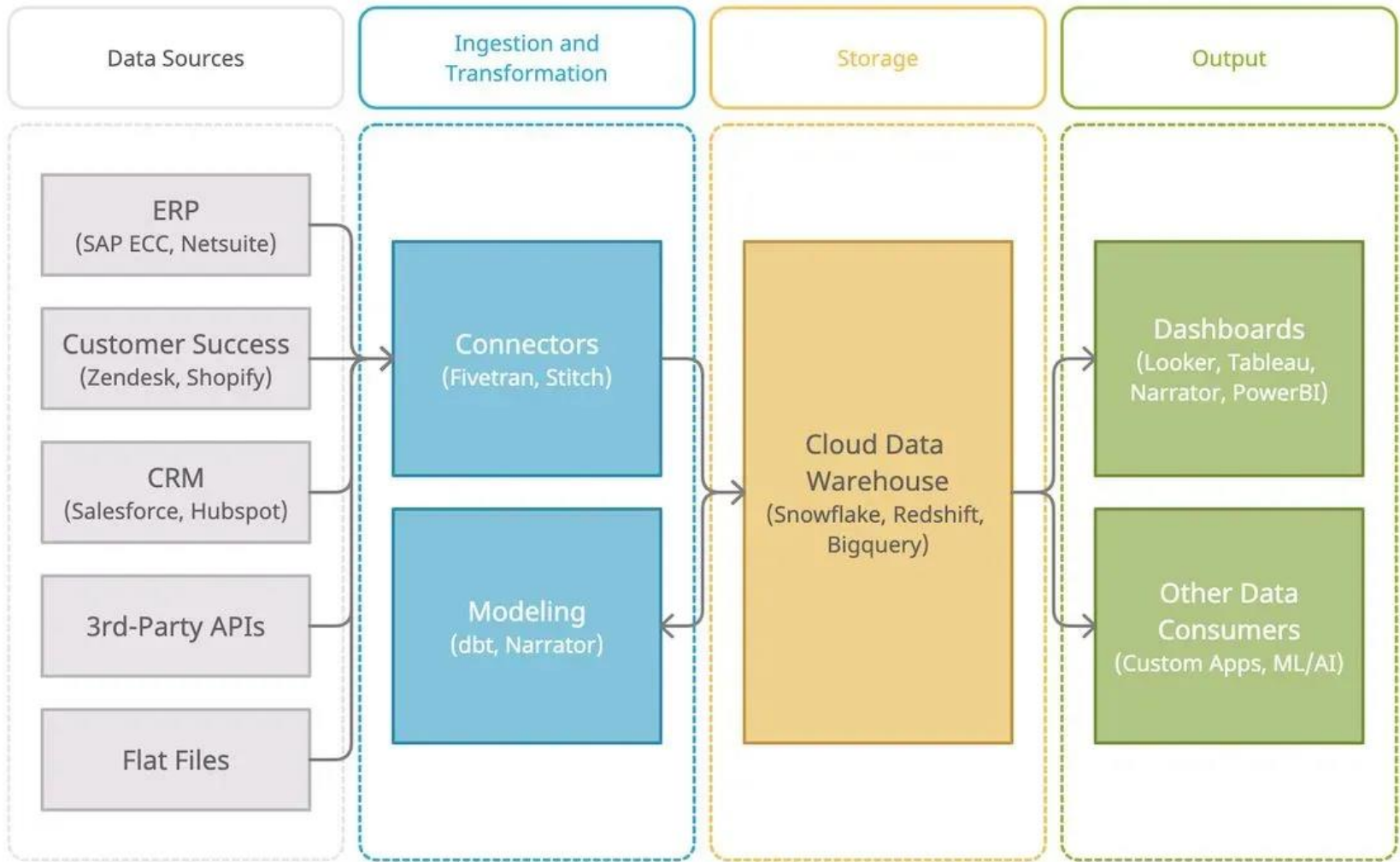
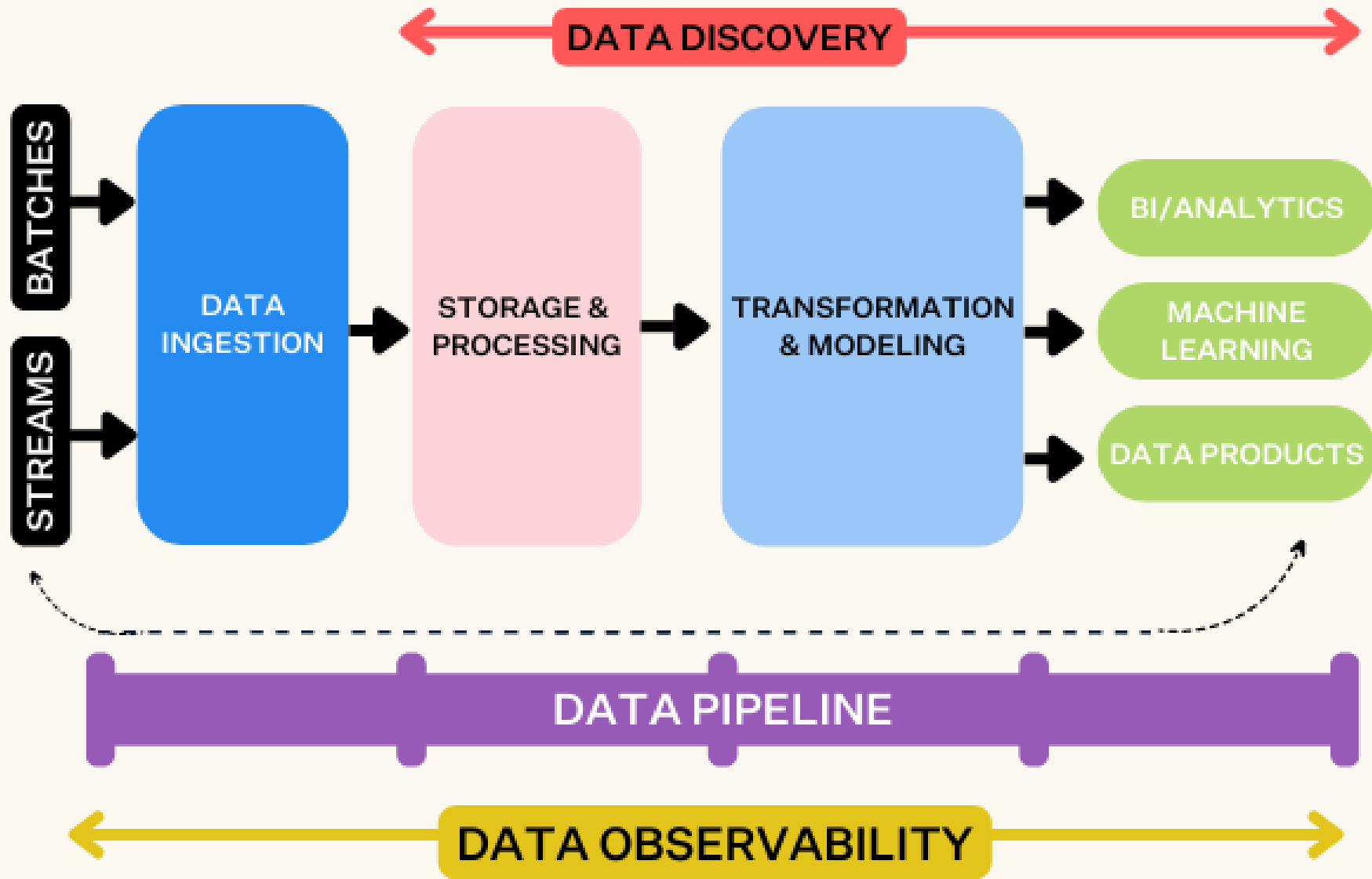
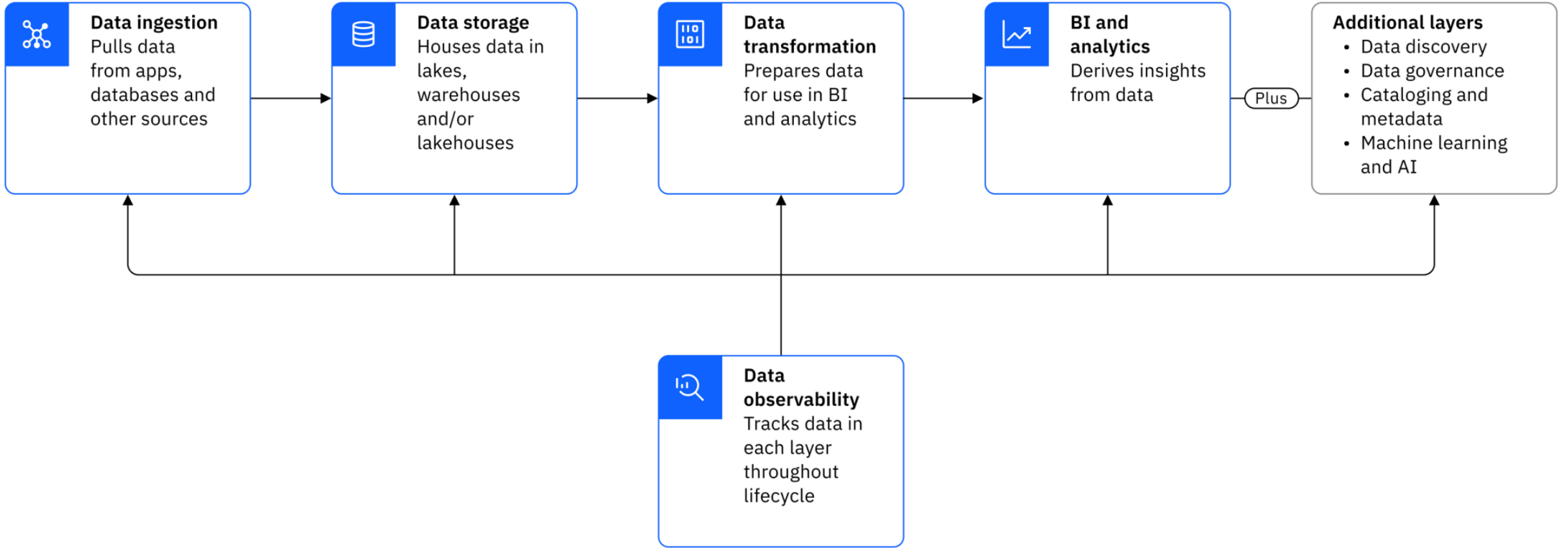
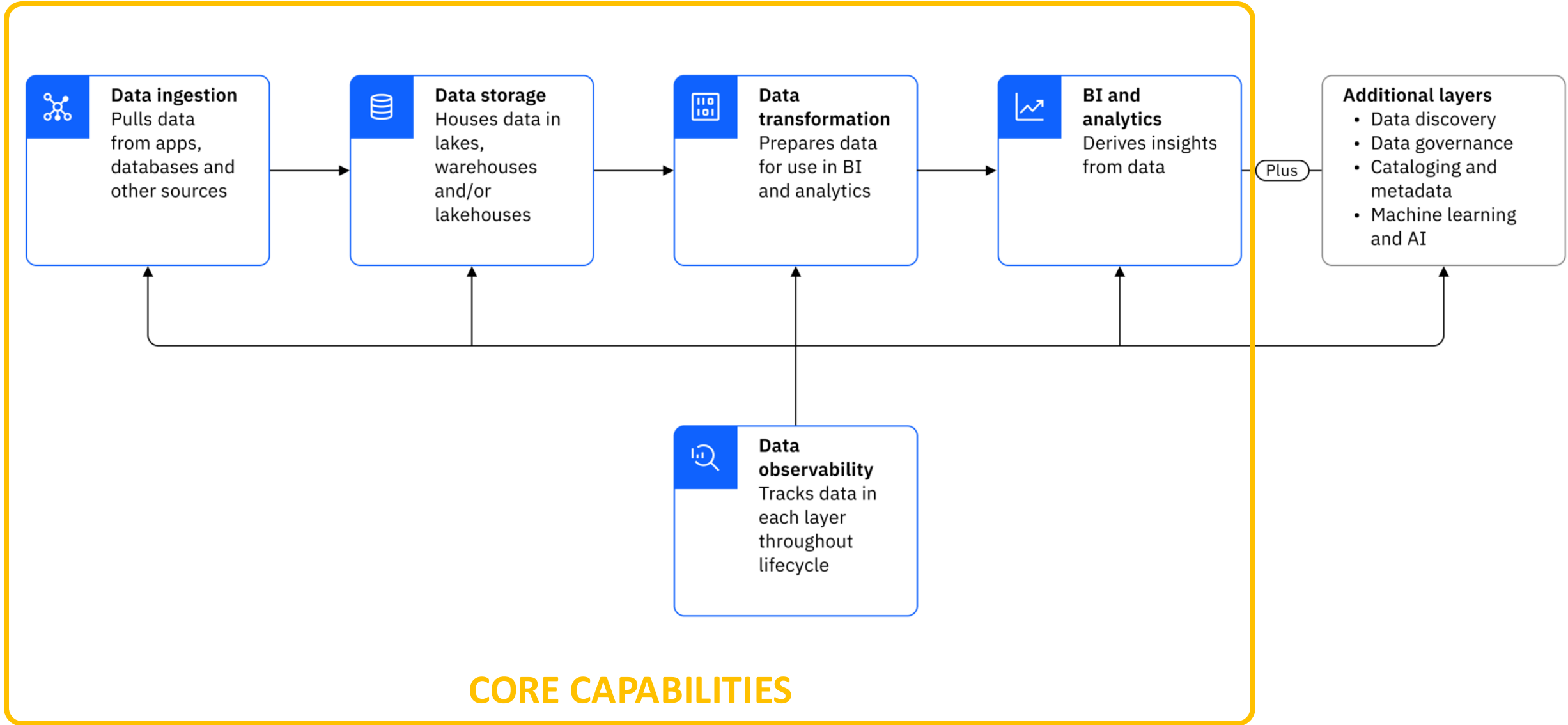


Illustration by Sanjiv Prasad



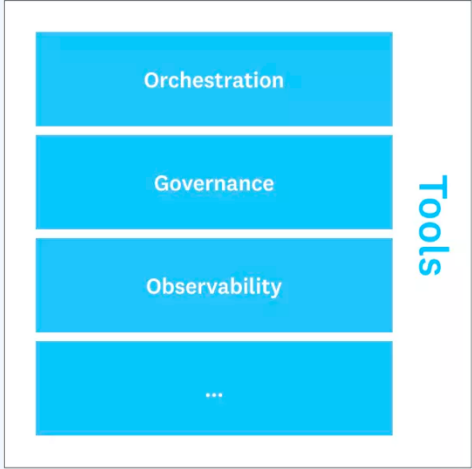
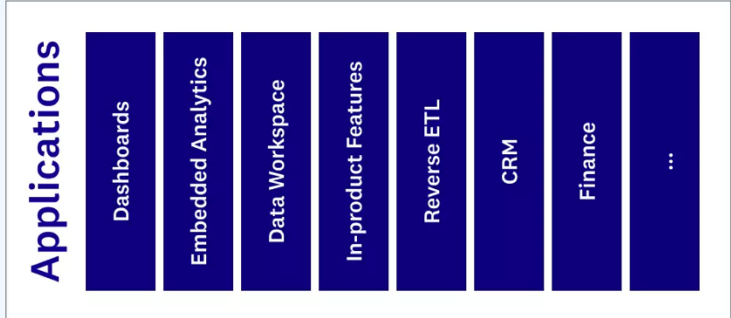
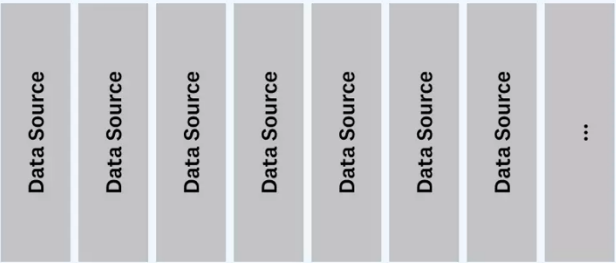




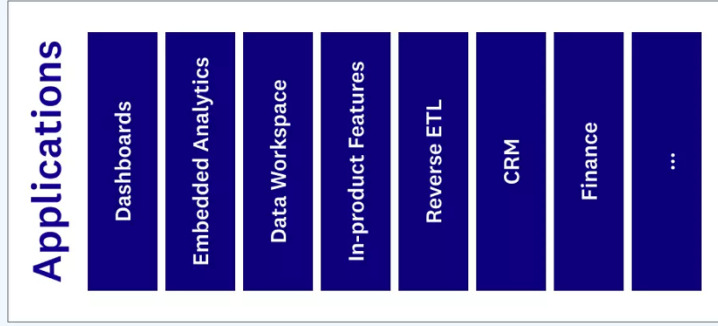


CORE CAPABILITIES

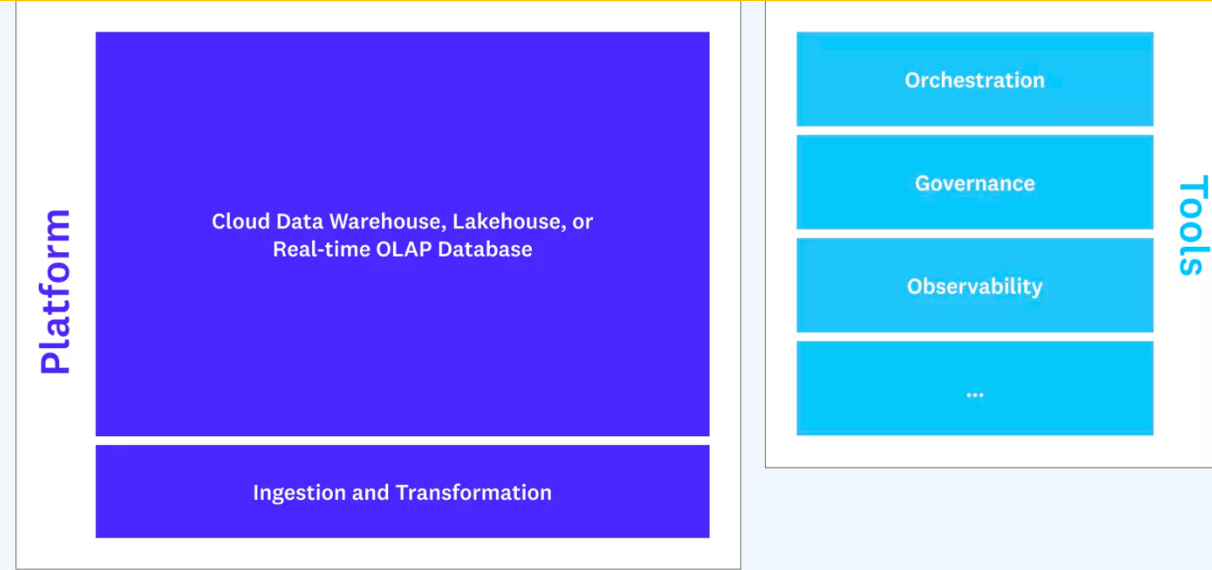




USE CASES



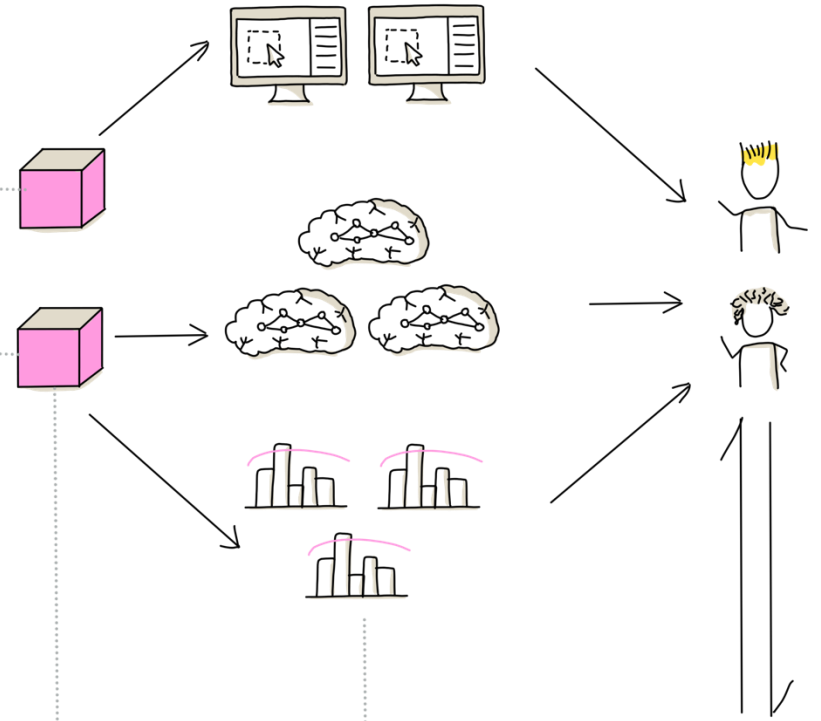
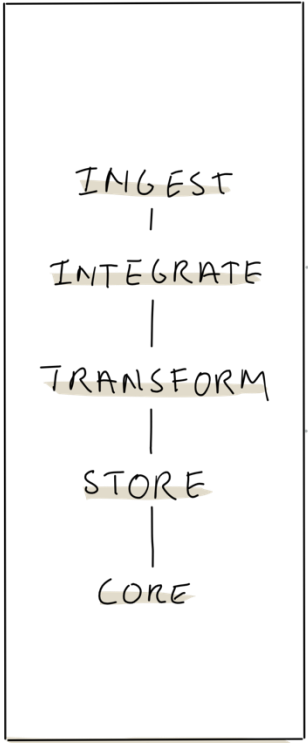
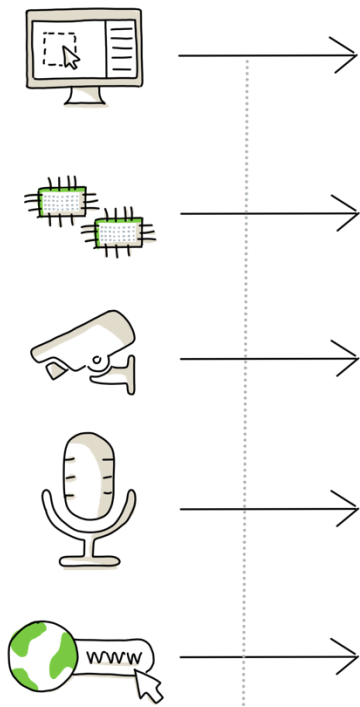
PLATFORM



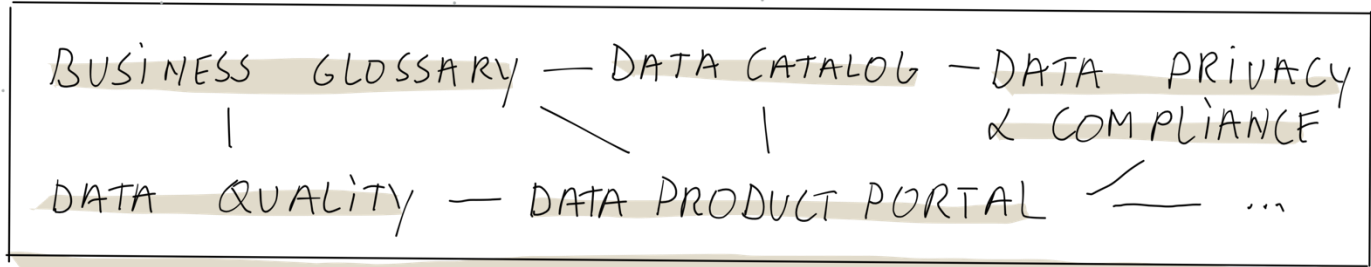
DATA PRODUCERS
(SOURCES)

DATA
PLATFORM

CONSUMPTION



VALUE



META-DATA MANAGEMENT

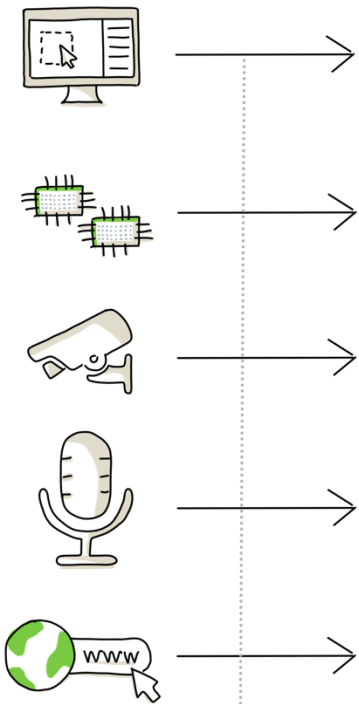


Table of Contents

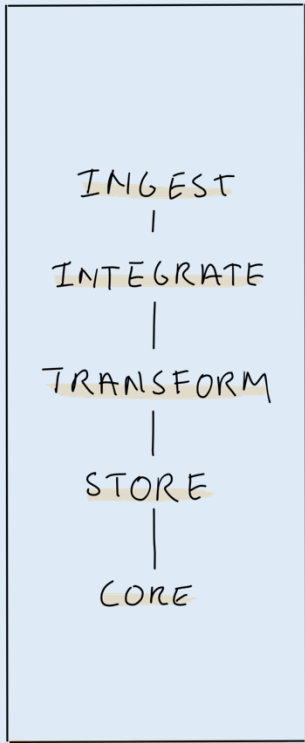
- Dead Horse Theory
- Data Platform
 - Introduction
 - **Core Layers**
 - Additional Layers
- Technology Selection



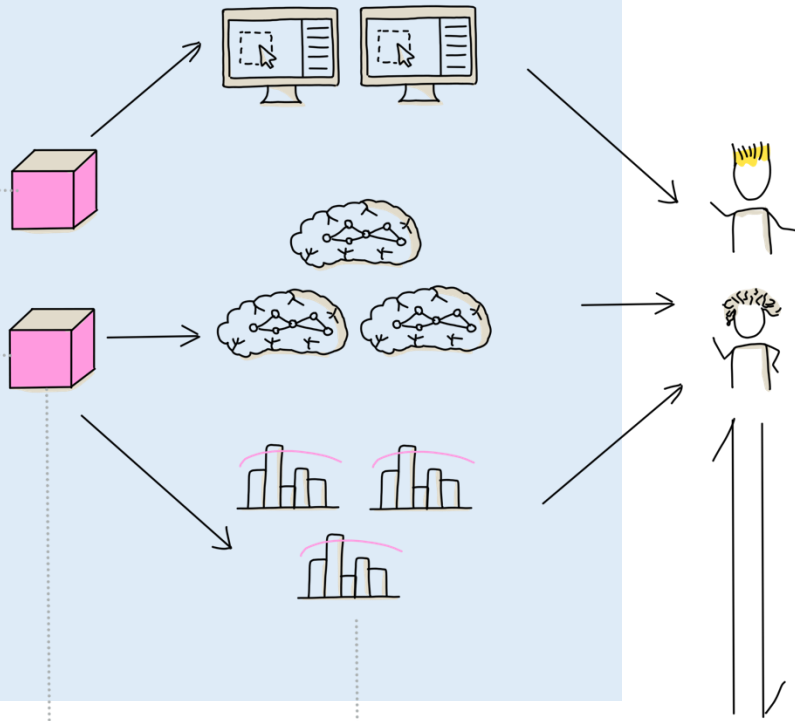
DATA PRODUCERS
(SOURCES)



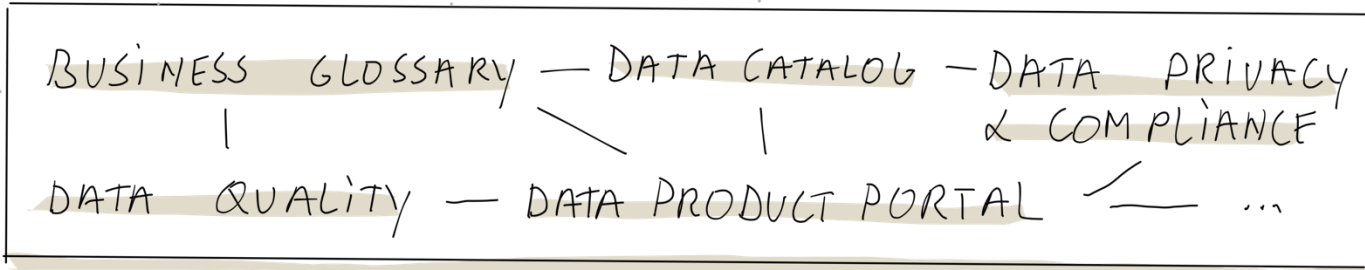
DATA PLATFORM



CONSUMPTION



VALUE



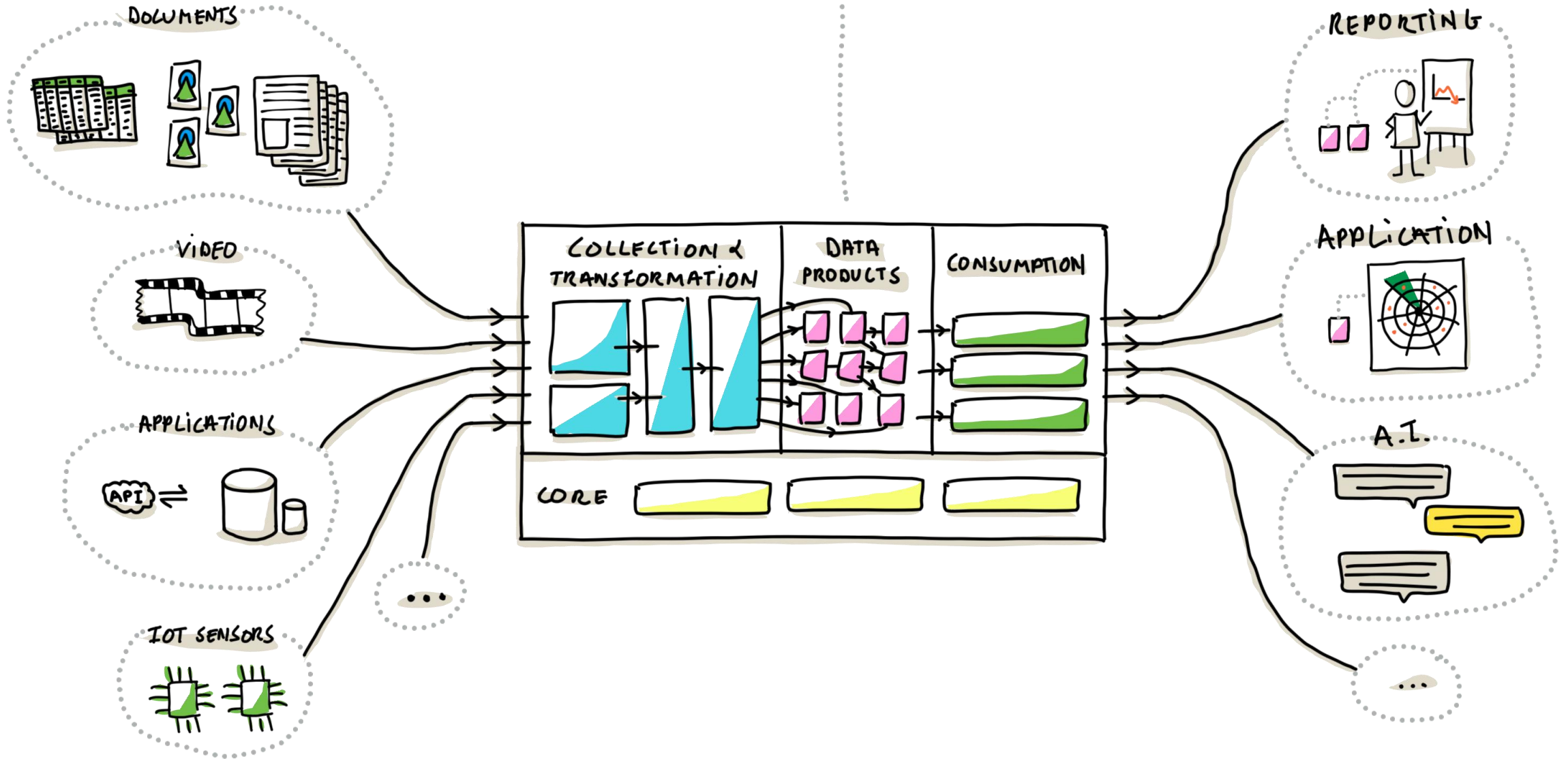
META-DATA MANAGEMENT



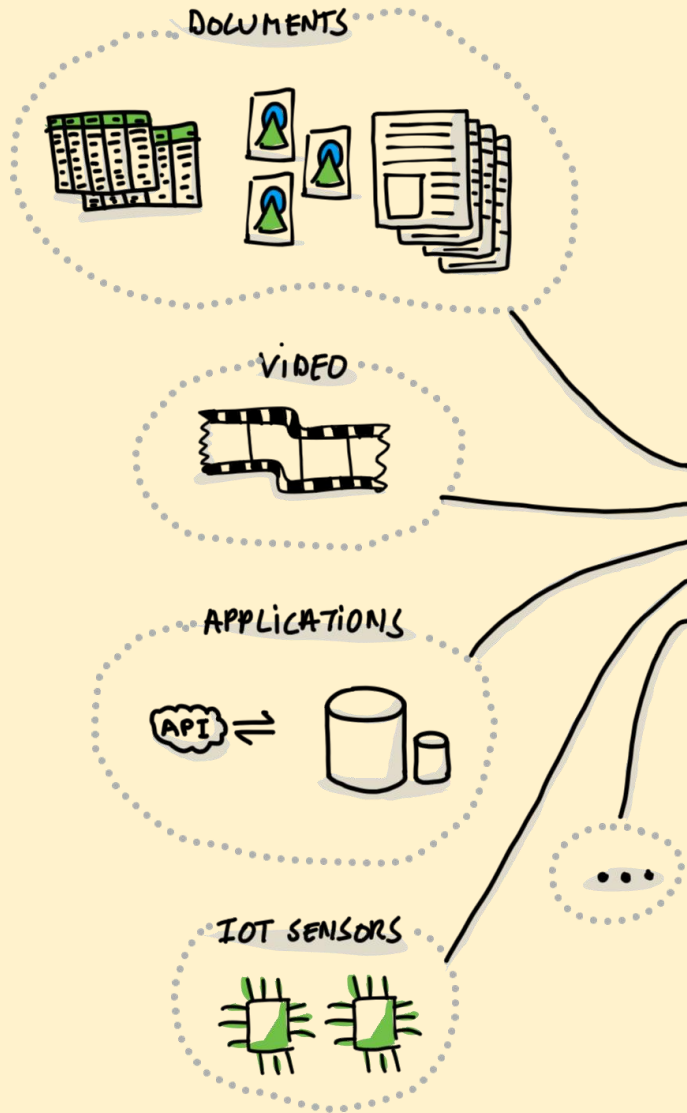
DATA SOURCES

DATA PLATFORM

CONSUMERS

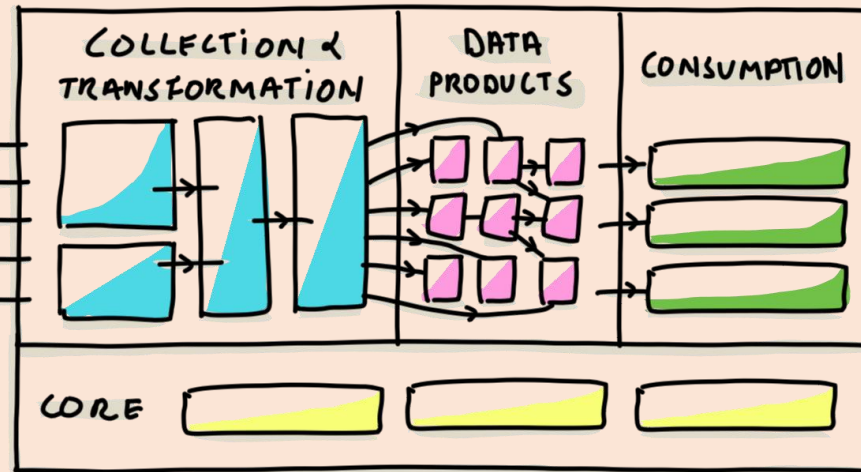


DATA SOURCES



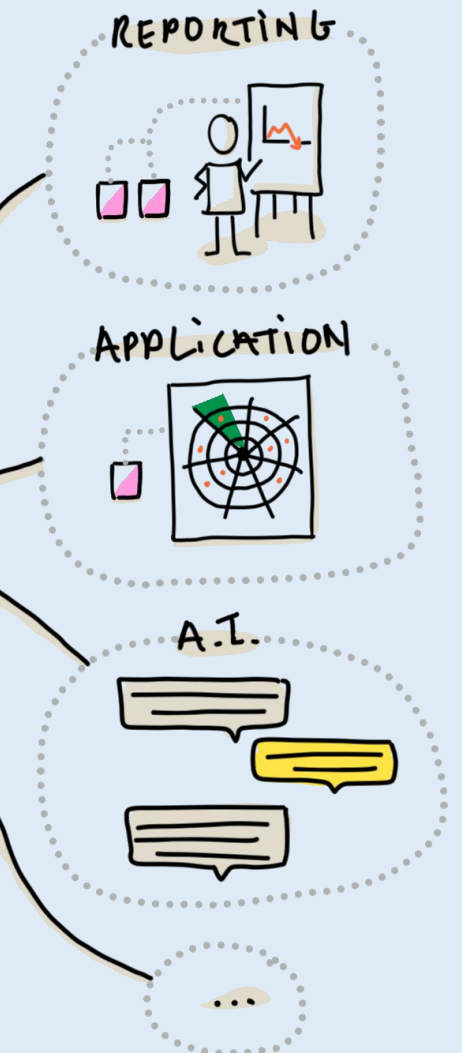
“Operational Plane”

DATA PLATFORM



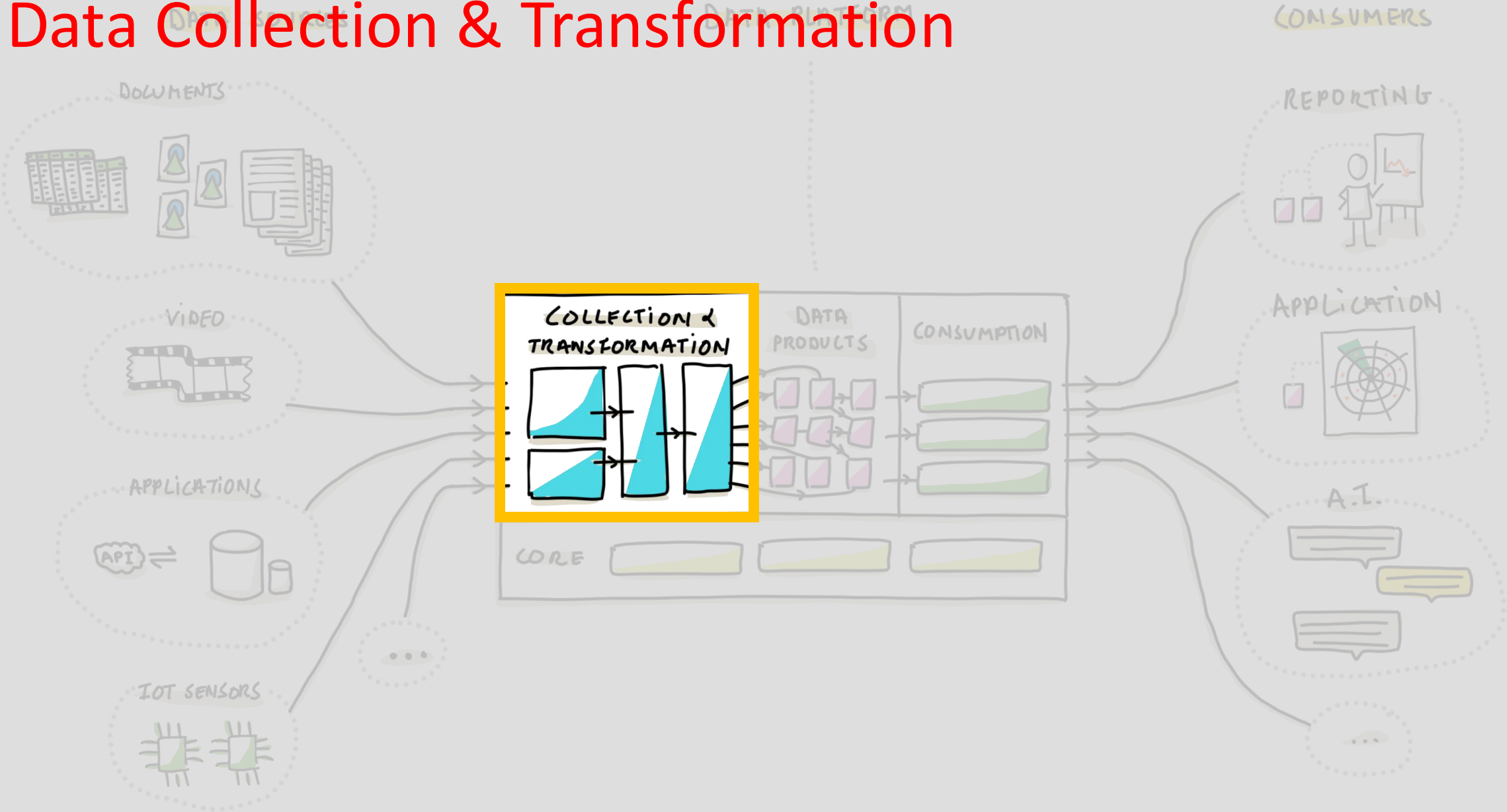
“Analytical Plane”

CONSUMERS



“Operational / Analytical Plane”

1. Data Collection & Transformation



“Operational Plane”

“Analytical Plane”

“Operational / Analytical Plane”

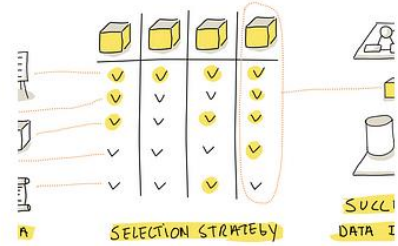


See also my Medium Blog

 janmeskens in The Modern Scientist


Data Ingestion—Part 2: Tool Selection Strategy

This article is the second one in my series on data ingestion. For an introduction to the topic and to explore 'data ingestion...



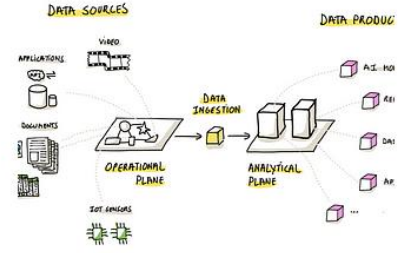
Jan 17  715  5



 janmeskens in The Modern Scientist

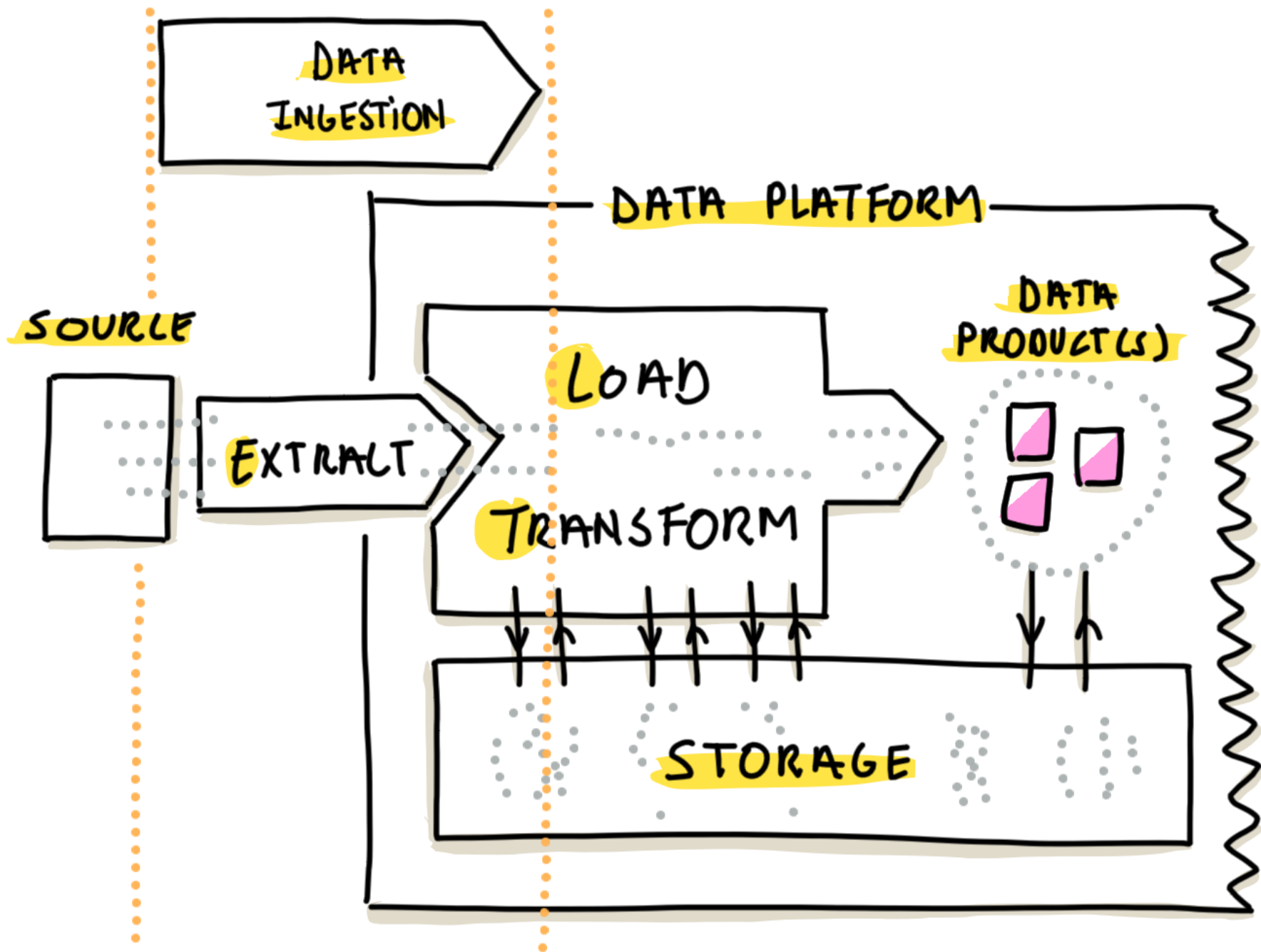
Data Ingestion — Part 1: Architectural Patterns

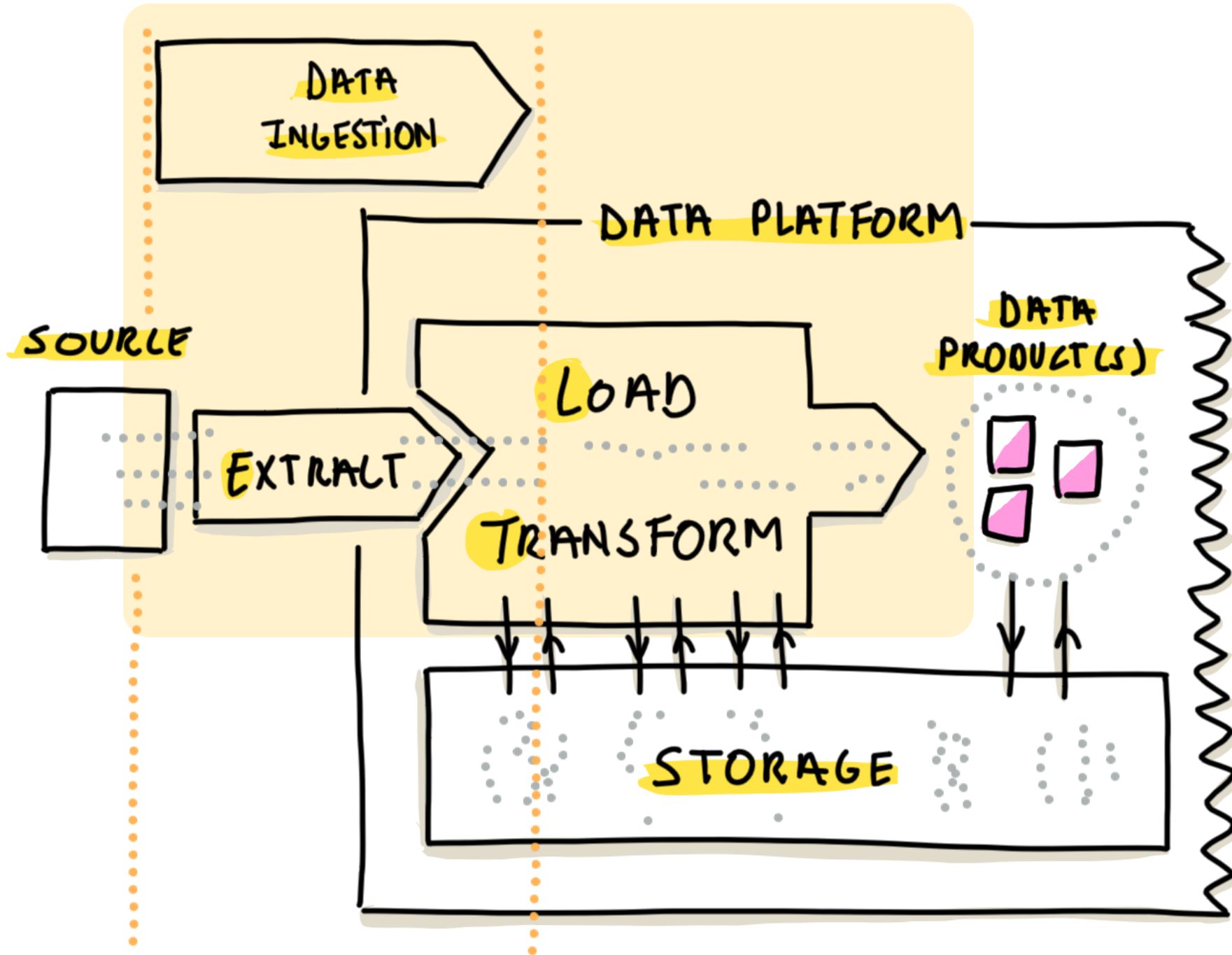
Over the course of two articles, I will thoroughly explore data ingestion, a fundamental process that bridges the operational...



Nov 27, 2023  2.7K  26

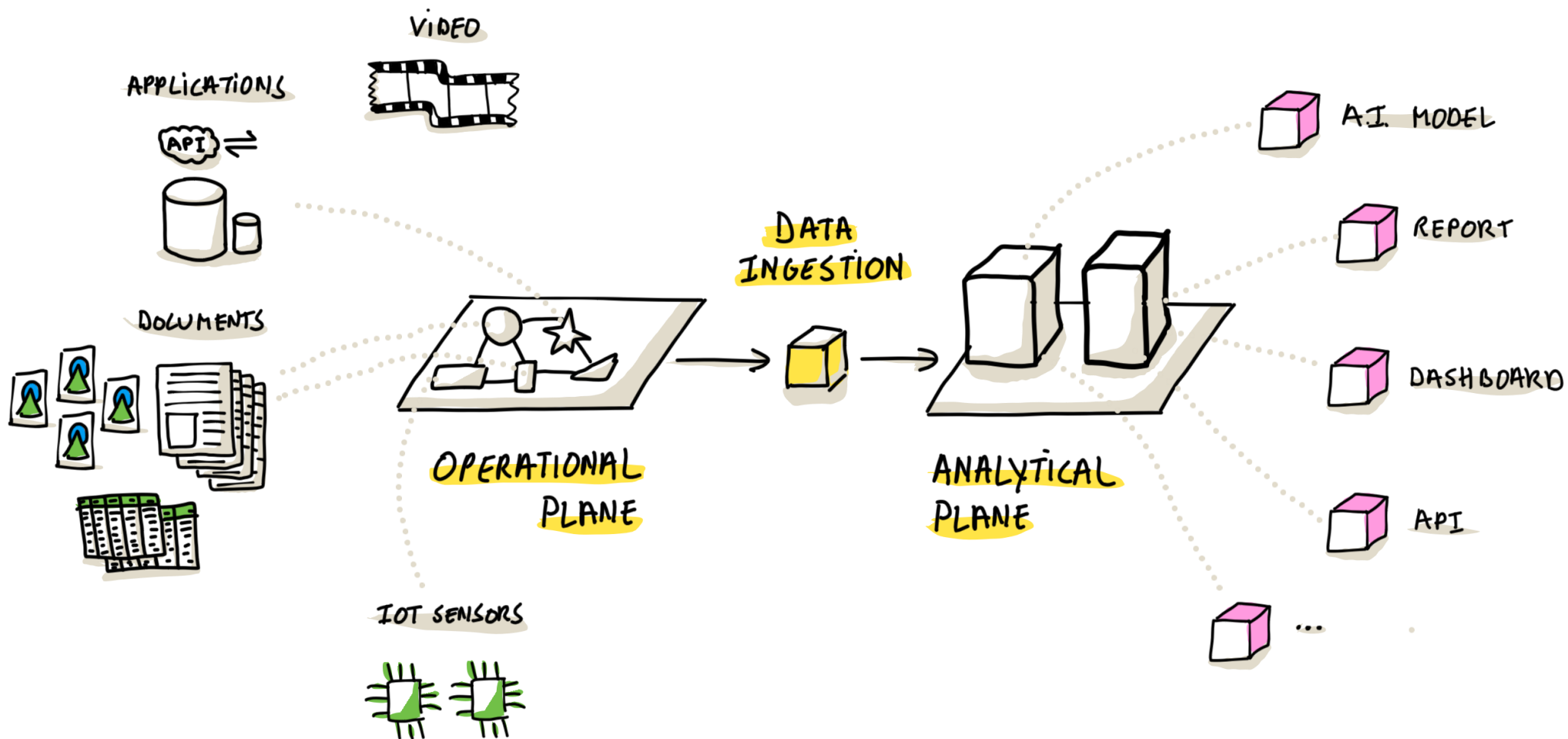






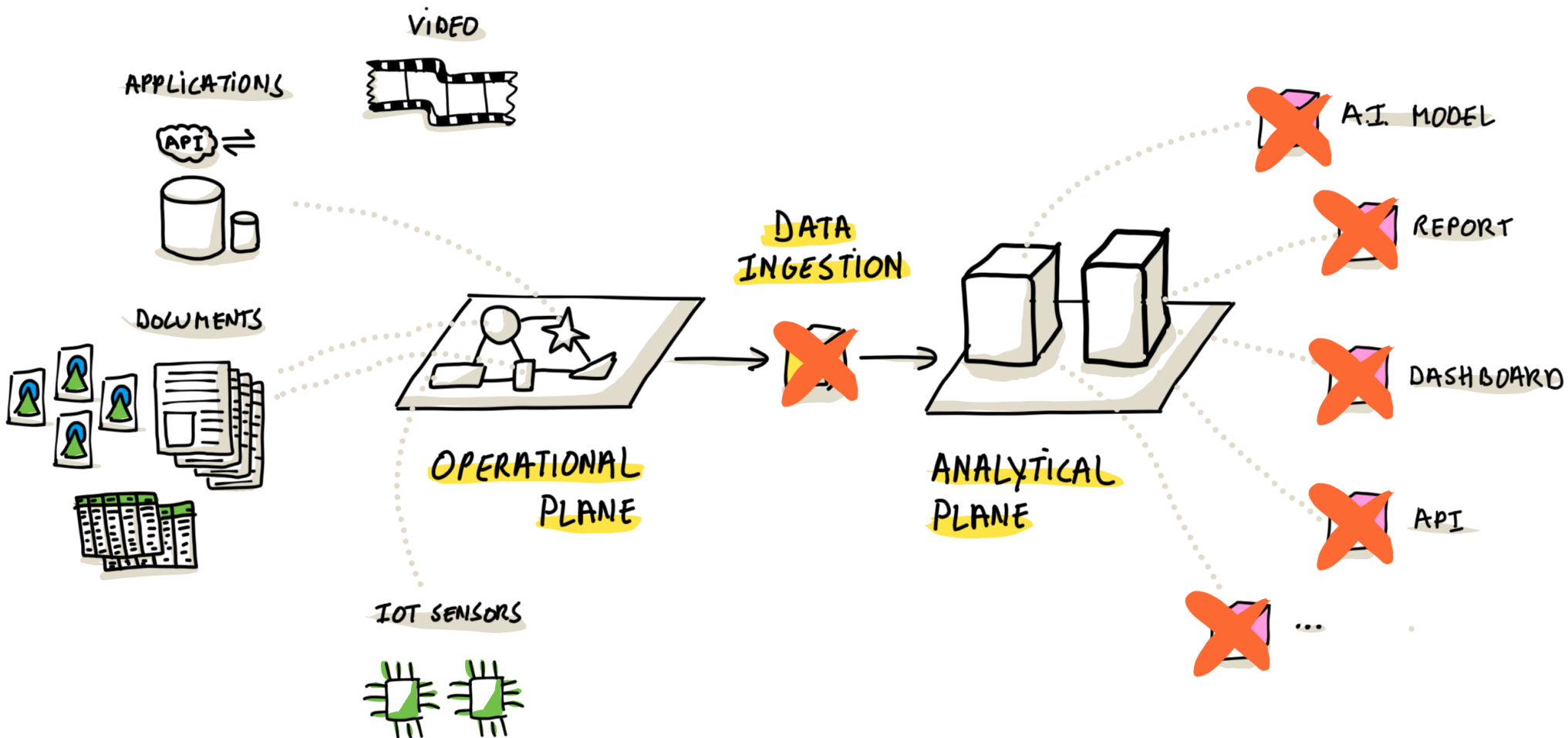
DATA SOURCES

DATA PRODUCTS



DATA SOURCES

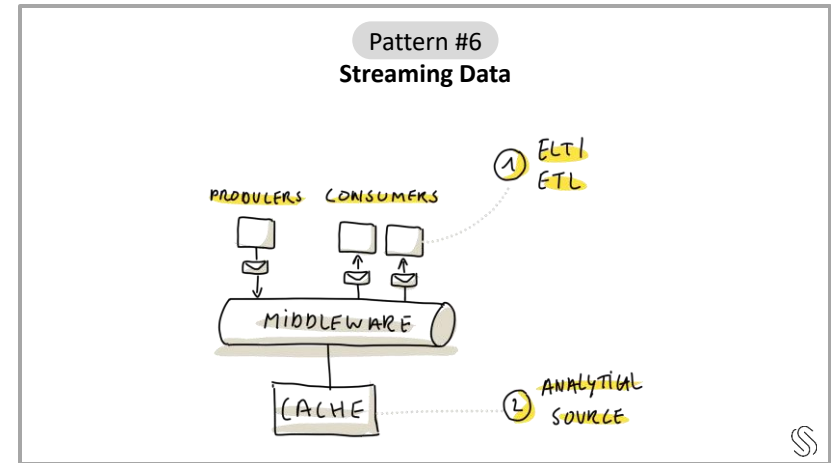
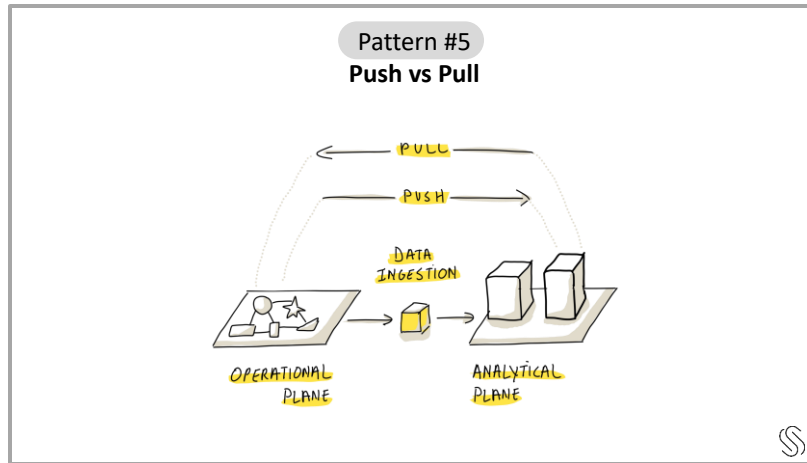
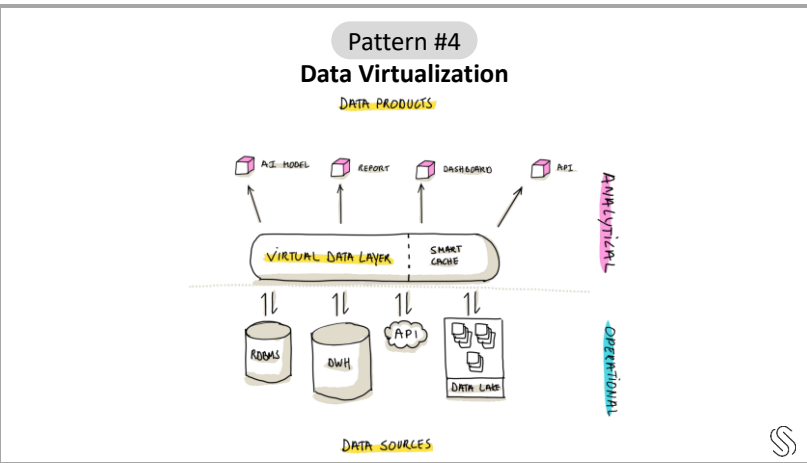
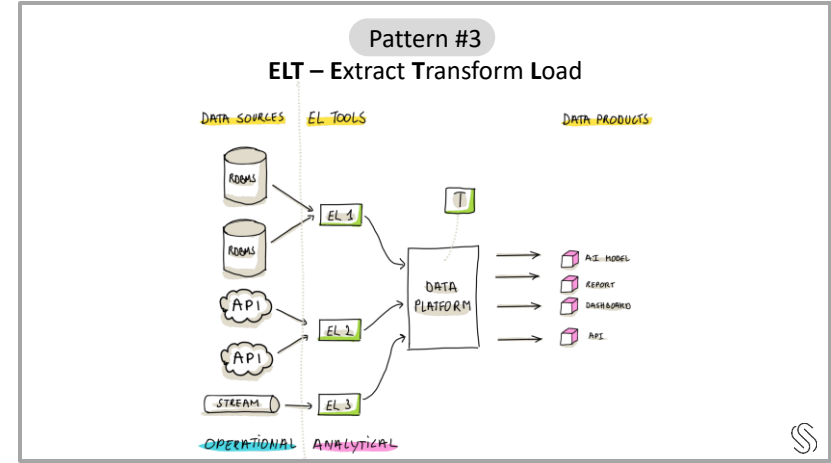
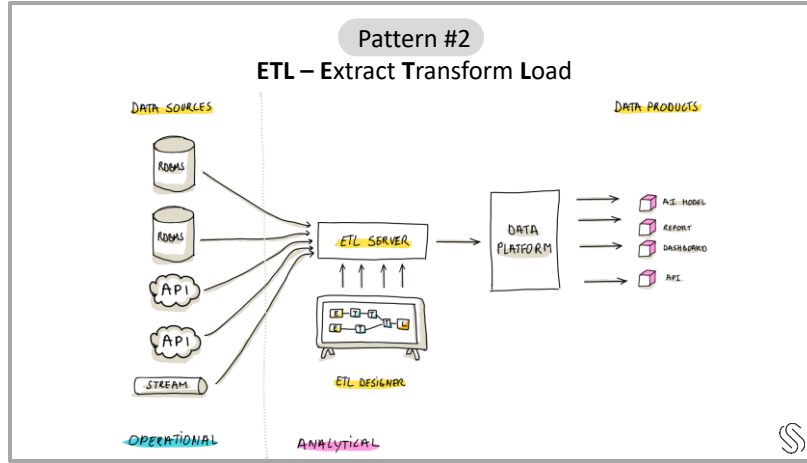
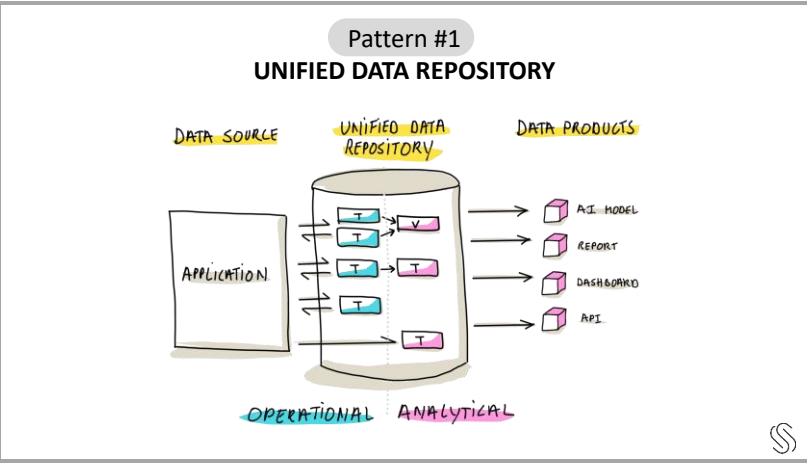
DATA PRODUCTS



Definition

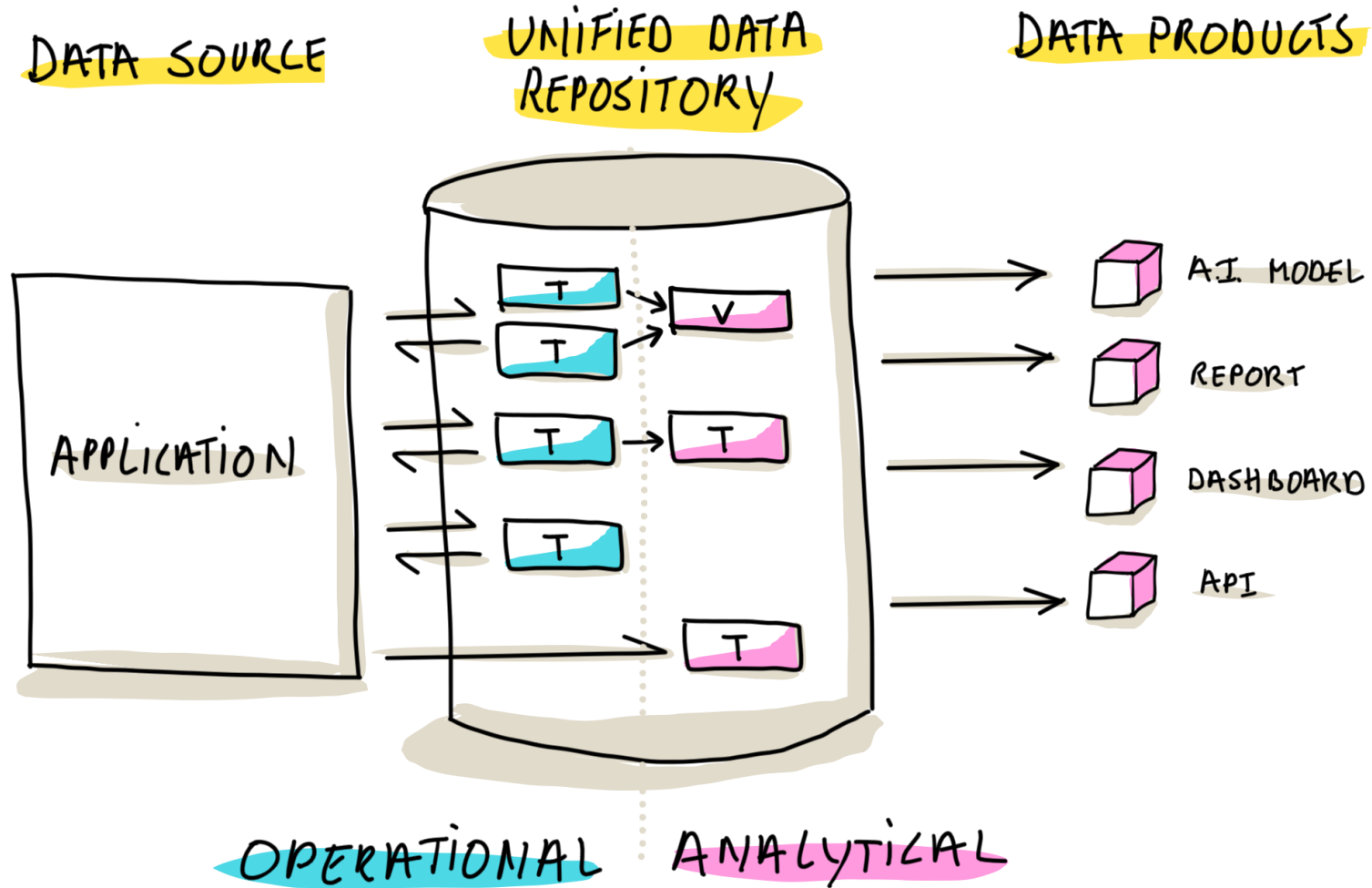
Data ingestion refers to the tools & processes used to **collect data from various sources and move it to a target site**, either in **batches or in real-time**. The data ingestion layer is **critical to your downstream** data science, BI, and analytics systems which depend on **timely, complete, and accurate data**.

Data Ingestion Patterns



Pattern #1

UNIFIED DATA REPOSITORY



Unified Data Repository = A single storage system caters to both the operational application needs and analytical processing

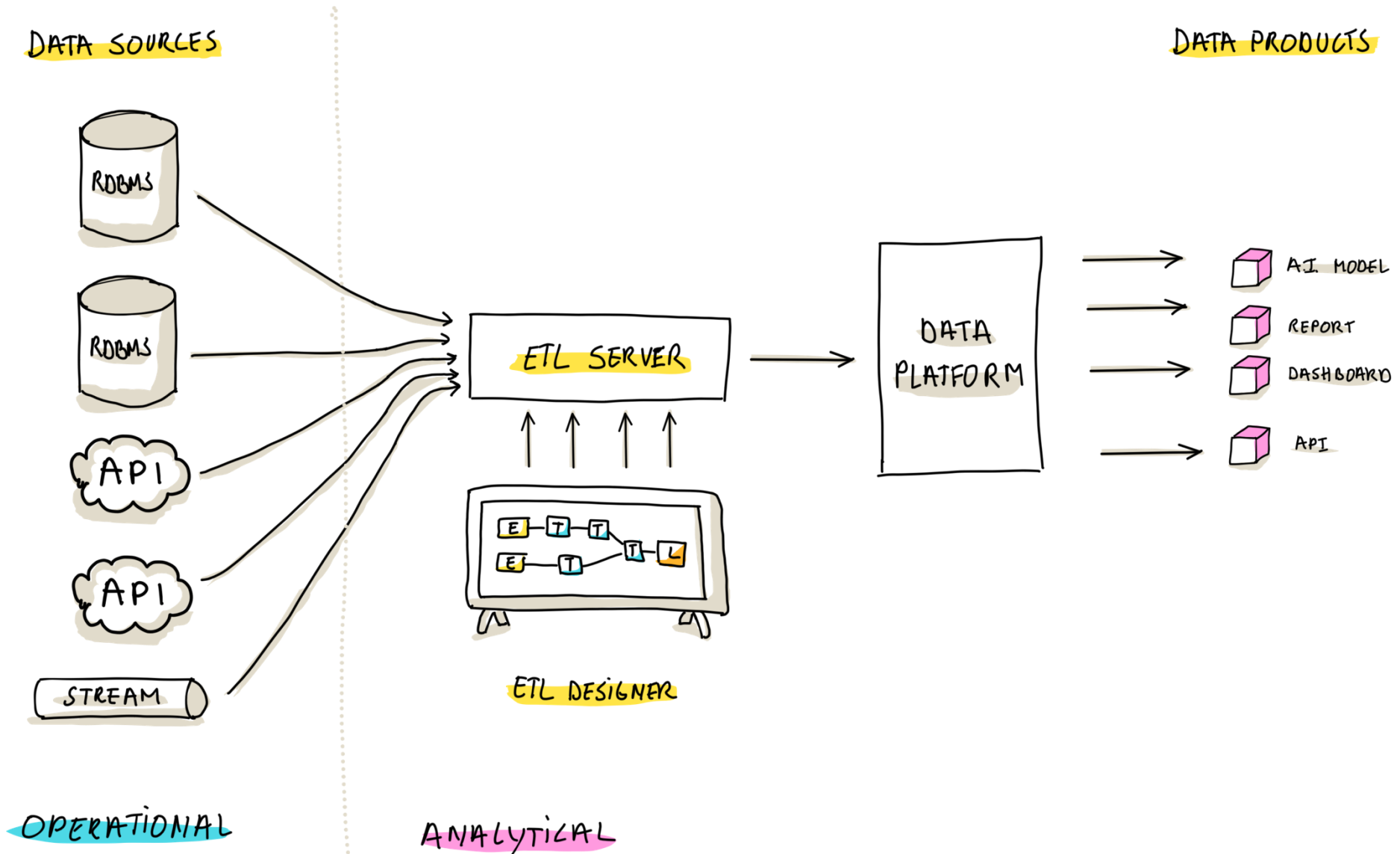
Specificities:

- Typically a Relational Database Management System (RDBMS).
- The same database is utilized for both everyday operations and data analysis
- Two prevalent sub-patterns:
 1. **Virtualization**
 2. **Duplication and Transformation**



Pattern #2

ETL – Extract Transform Load



ETL (Extract, Transform, Load) = *a well-established paradigm in data processing*

3 steps:

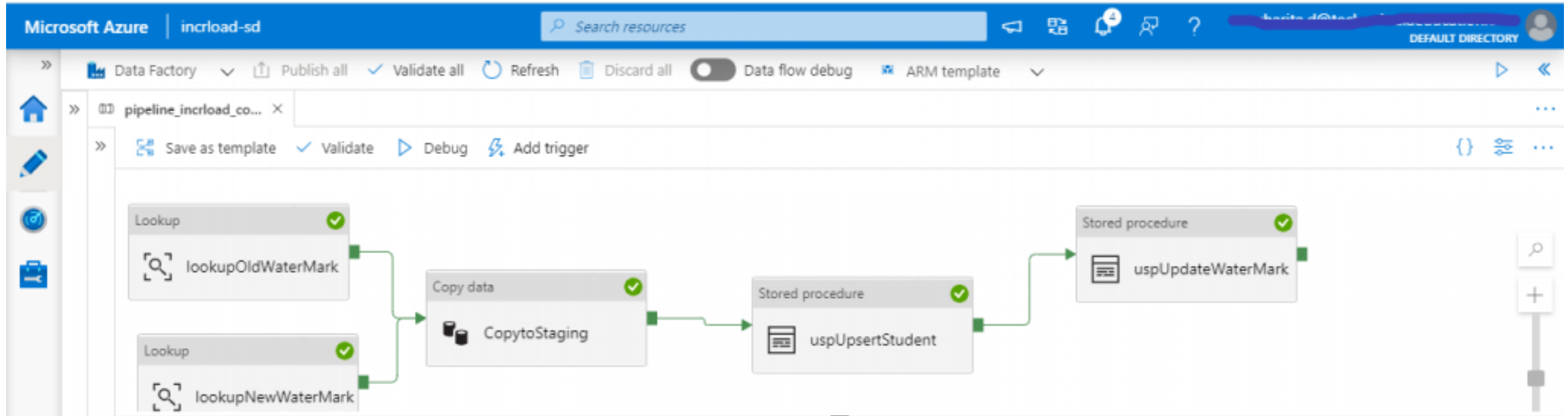
1. **Extract** : Data is harvested from its source
2. **Transform** : Refined on an ETL server
3. **Load** : The polished output is deposited into an analytics-focused database.

ETL tools have a graphical interface where users can interlink Extract, Transform, and Load operations within an intuitive visual workflow. These processes are often further customizable through scripting or direct SQL queries.



Pattern #2

ETL – Extract Transform Load



Extract Transform

Load

Pattern #2

ETL – Extract Transform Load

Microsoft Azure | incload-sd

Data Factory | Publish all | Validate

pipeline_incload_co...

Save as template | Validate

Lookup
lookupOldWaterMark

Lookup
lookupNewWaterMark

Copy data
CopytoStaging

Stored procedure
uspUpsertStudent

Stored procedure
uspUpdateWaterMark

General Settings User properties

Source dataset * SqlServerTable1 Open + New Preview data

Use query Table Query Stored procedure

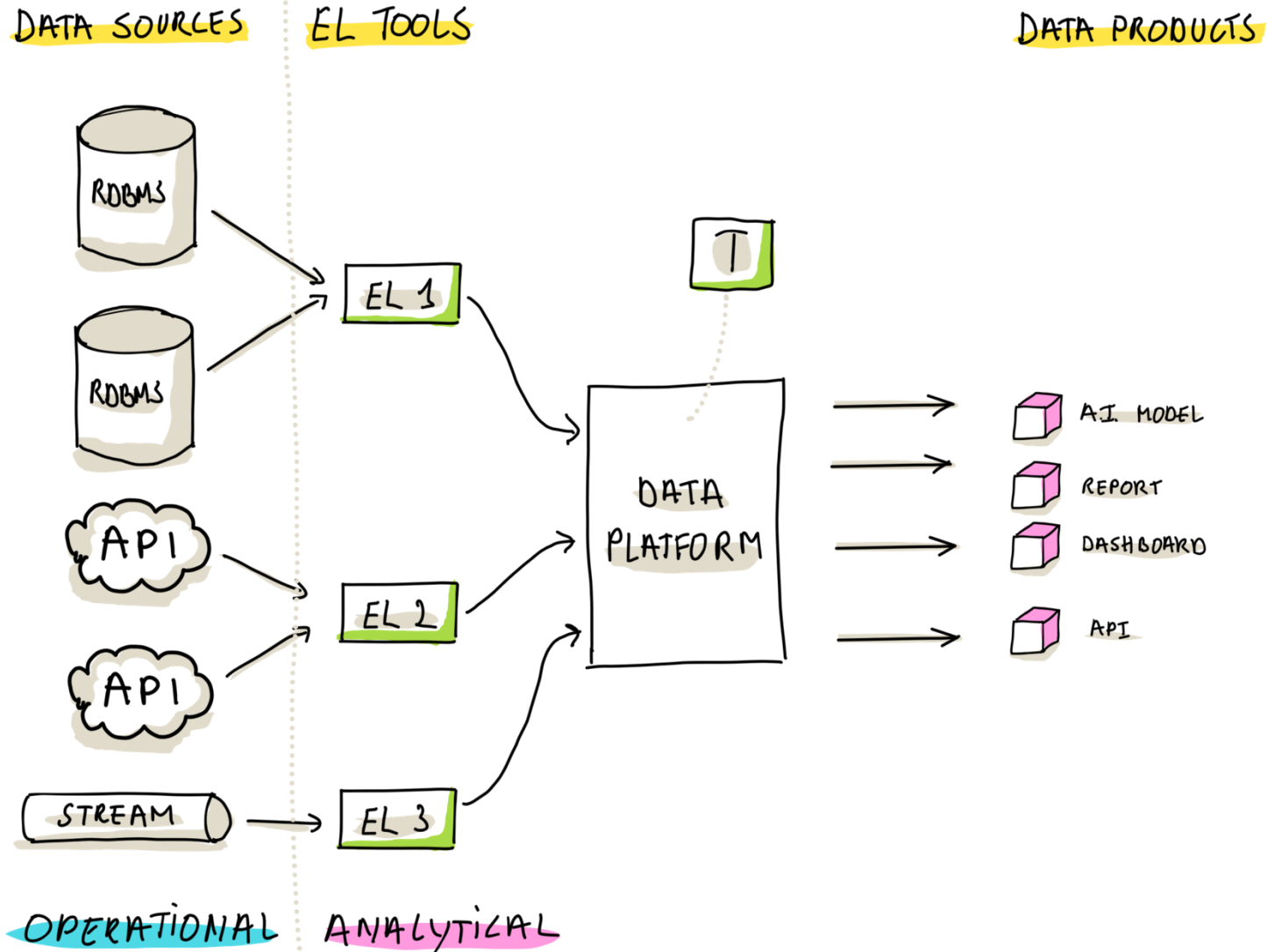
Query
SELECT
MAX(@{pipeline().parameters.waterMarkCol}) AS NewwaterMarkVal FROM
@{pipeline().parameters.srcTableName}

Query timeout (minutes) 120

Extract Transform **Load**

Pattern #3

ELT – Extract Load Transform



ELT, sharing the basic steps of ETL, diverges by restructuring and redefining these processes.

In ELT:

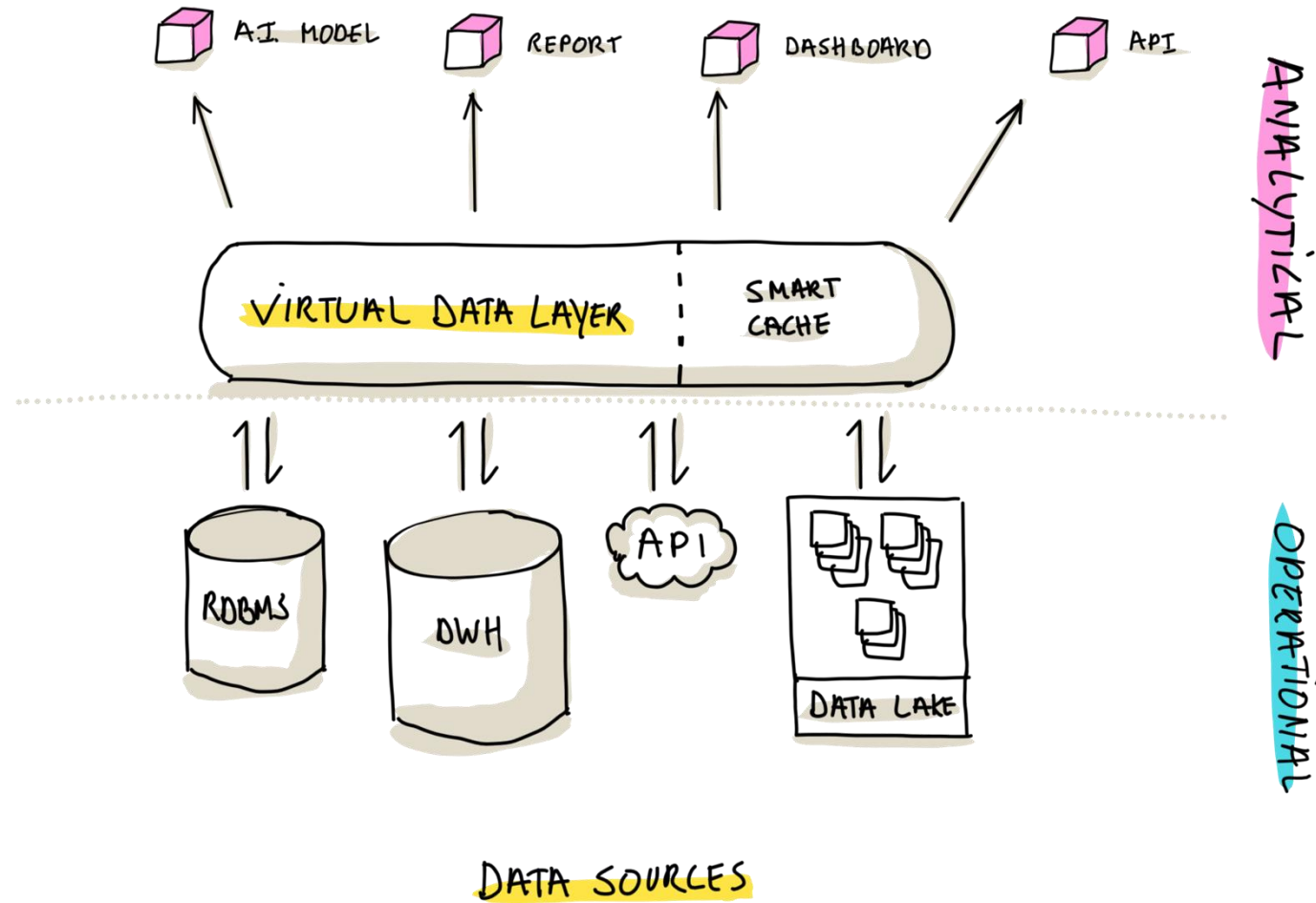
- 1. EL — Extract and Load operations** are carried out first, transferring raw data directly to the data platform without immediate transformation.
- 2. T — Transformation** occurs subsequently, converting raw data into actionable insights. Crucially, transformation tasks can operate independently and on different schedules from the extraction and loading.



Pattern #4

Data Virtualization

DATA PRODUCTS



Data Virtualization = specialized software to establish a virtualized data layer over multiple underlying data sources. This intermediary layer allows for the execution of queries that are partially processed by the original data sources, integrating the results into a cohesive dataset for analysis.

- Inspired by the Unified Data Repository (Pattern #1)



BI integration

Client support

Starburst Enterprise

Analytics engine

MPP query engine Data products Fault-tolerant execution

Query optimizer Elastic auto scaling Smart indexing & caching Metrics & logging

Global security

Fine-grained access control End-to-end encryption Data masking Event logging Query auditing

Data lakes

Relational DBs

NoSQL stores

Real-time analytics

Applications

Any data source anywhere
Cross-cloud and region analytics

On Premise

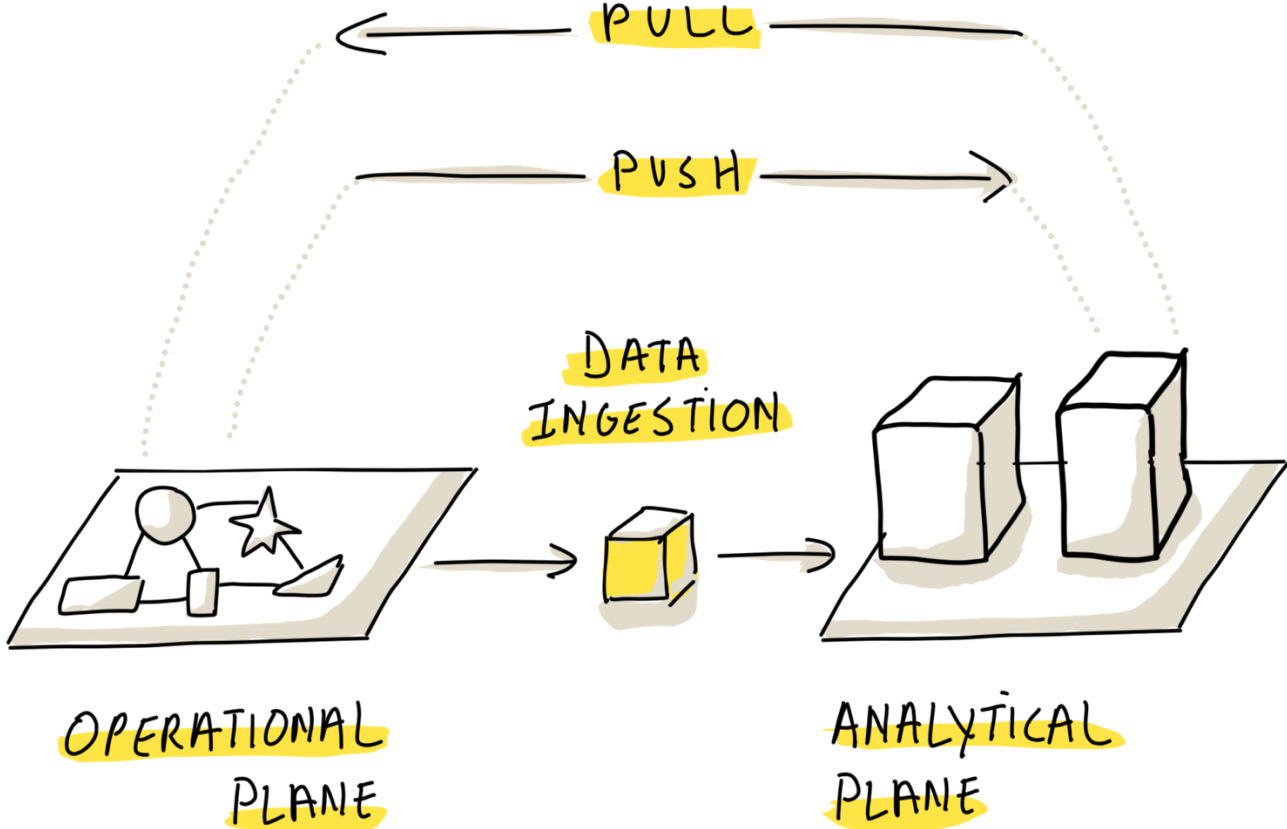
Hybrid

Cross-cloud



Pattern #5

Push vs Pull



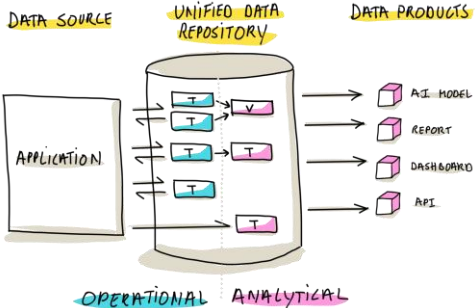
Push = The operational plane initiates data transfer to an endpoint designated by the analytical plane.

- Often found within streaming architectures (discussed next) but is not confined to them.
- Software development teams are mostly responsible to implement the push mechanism



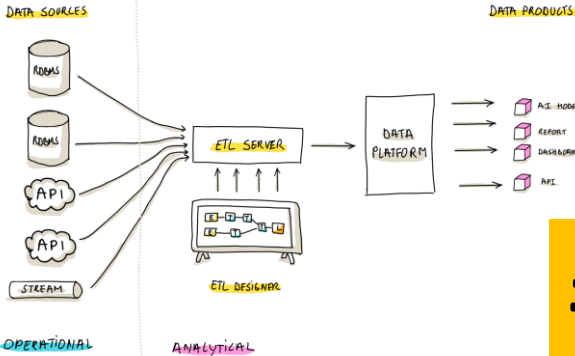
EXERCISE : STRENGTHS & WEAKNESSES

Pattern #1
UNIFIED DATA REPOSITORY



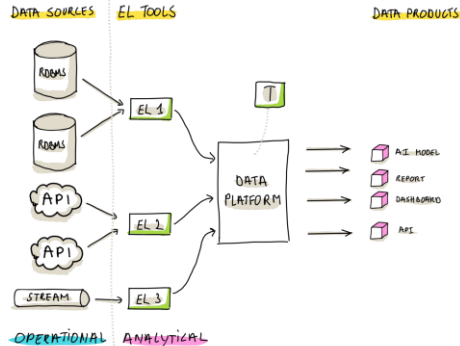
#1

Pattern #2
ETL – Extract Transform Load



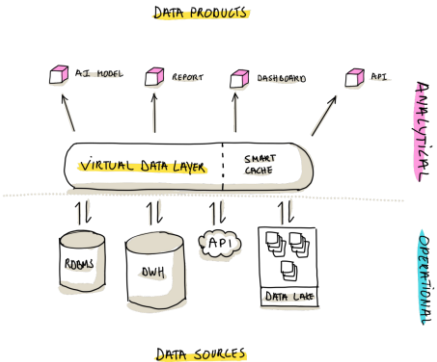
#2

Pattern #3
ELT – Extract Transform Load



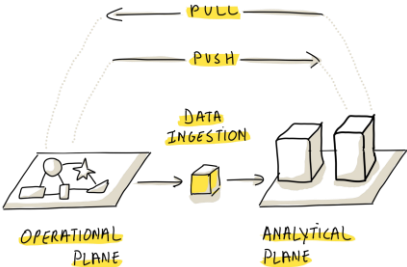
#3

Pattern #4
Data Virtualization



#4

Pattern #5
Push vs Pull



#5



Pattern #6

Streaming Data

**Streams
record history**



“The sequence of moves”

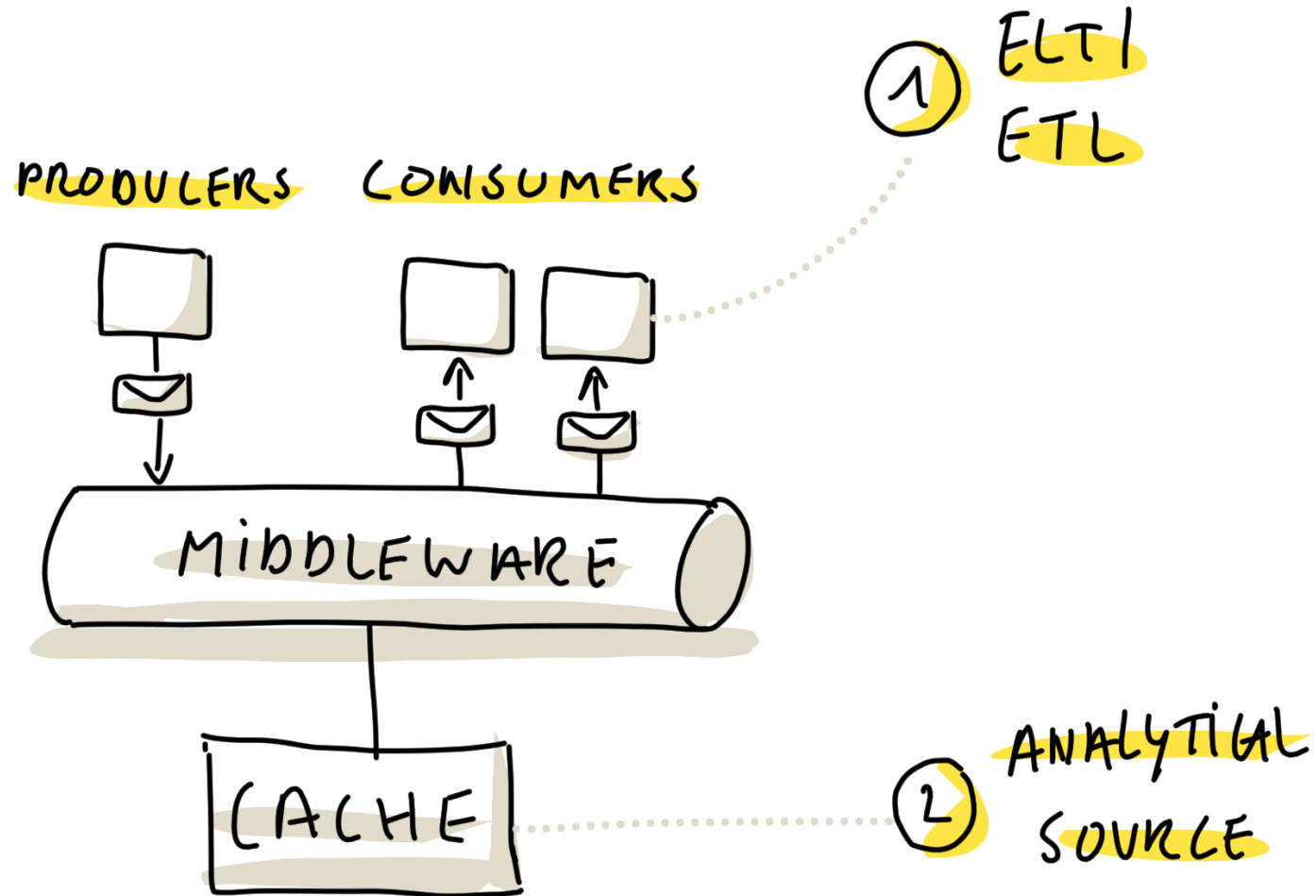
**Tables
represent state**



“The state of the board”



Pattern #6 Streaming Data



Stream processing = the continuous flow of data as it's generated, enabling real-time processing and analysis for immediate insights.

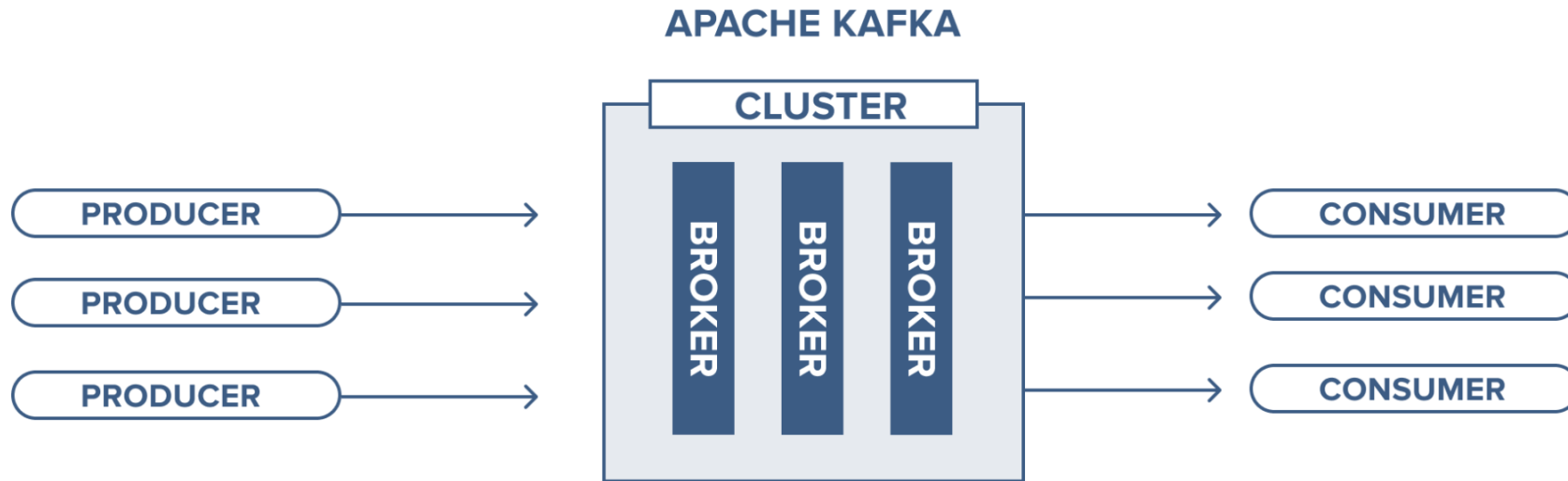
These systems are crucial for instant decision-making tasks and support high-volume, low-latency processing for activities like financial trades, real-time analytics, and IoT monitoring.

Two common approaches:

- ELT (or ETL) for streaming
- Leveraging streaming caches

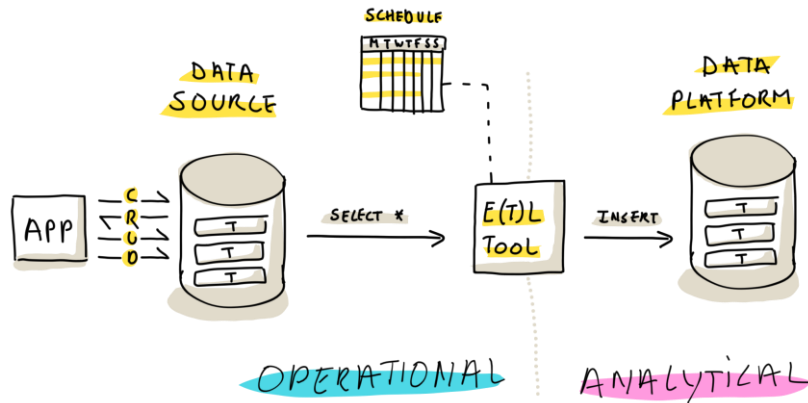


Example: Kafka

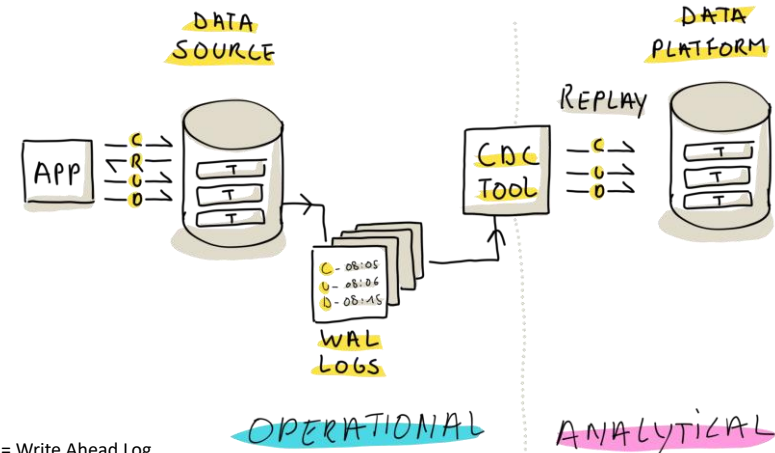


Data Ingestion Tool Flavors

Flavor #1
Batch Loading



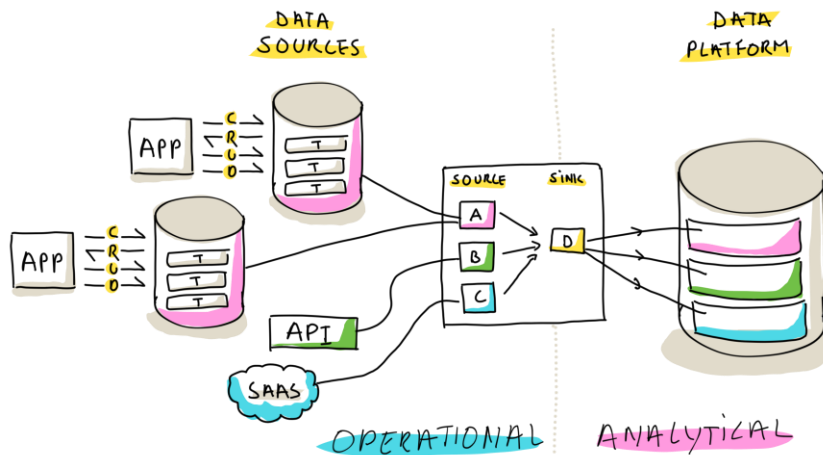
Flavor #2
CDC – Change Data Capture



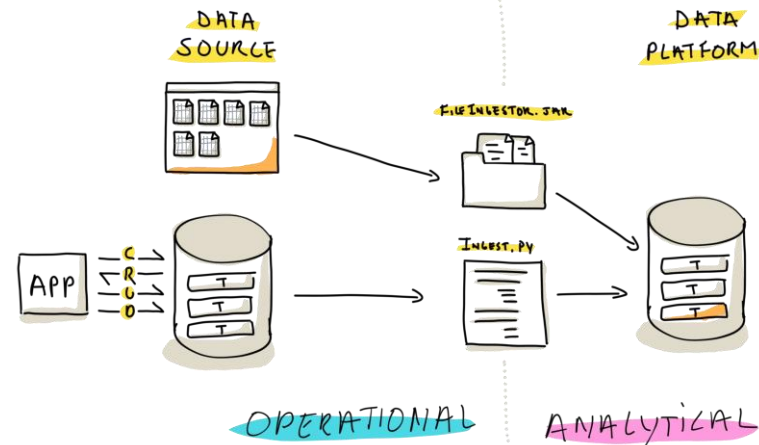
WAL = Write Ahead Log



Flavor #3
Connector Based

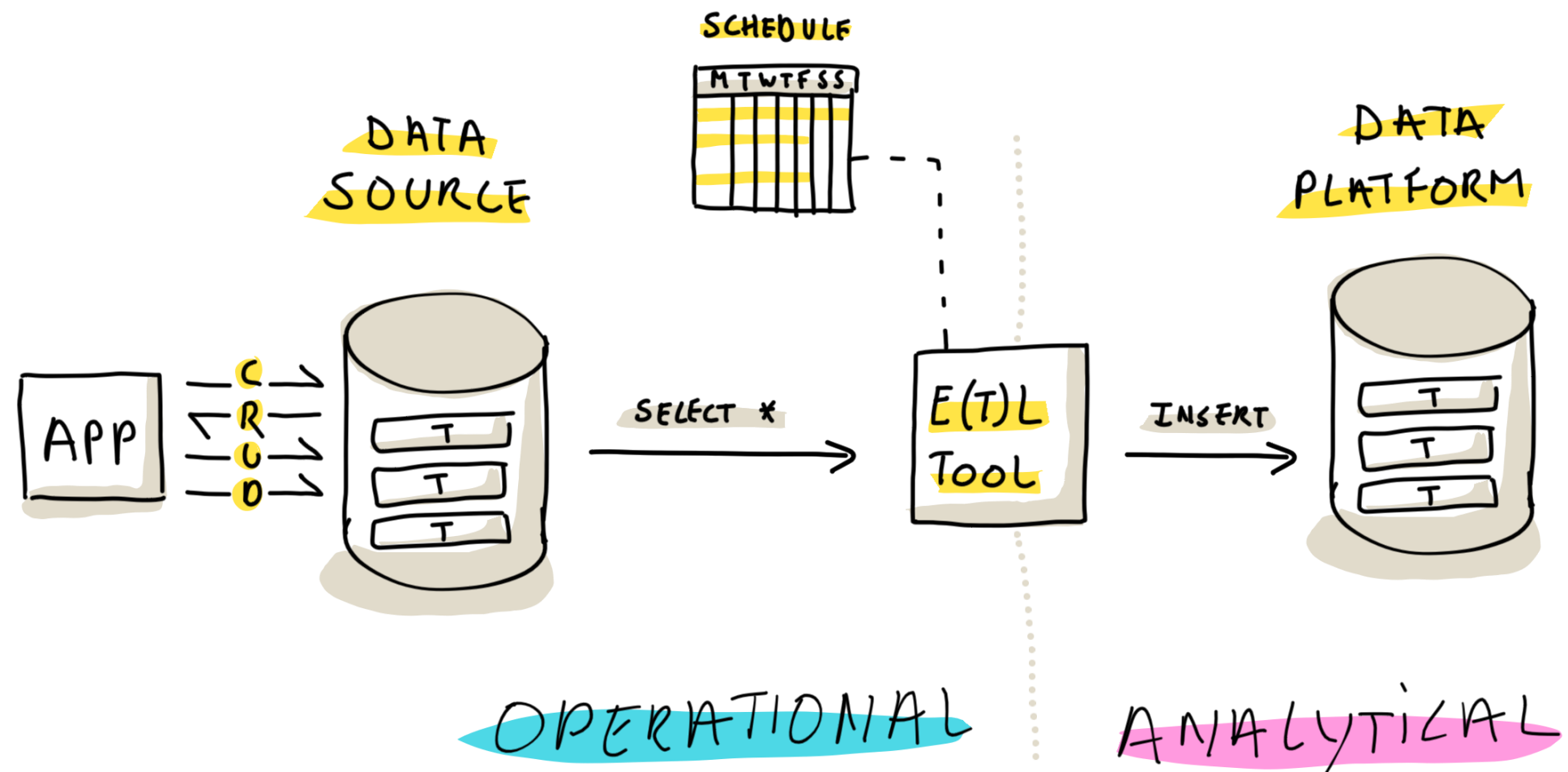


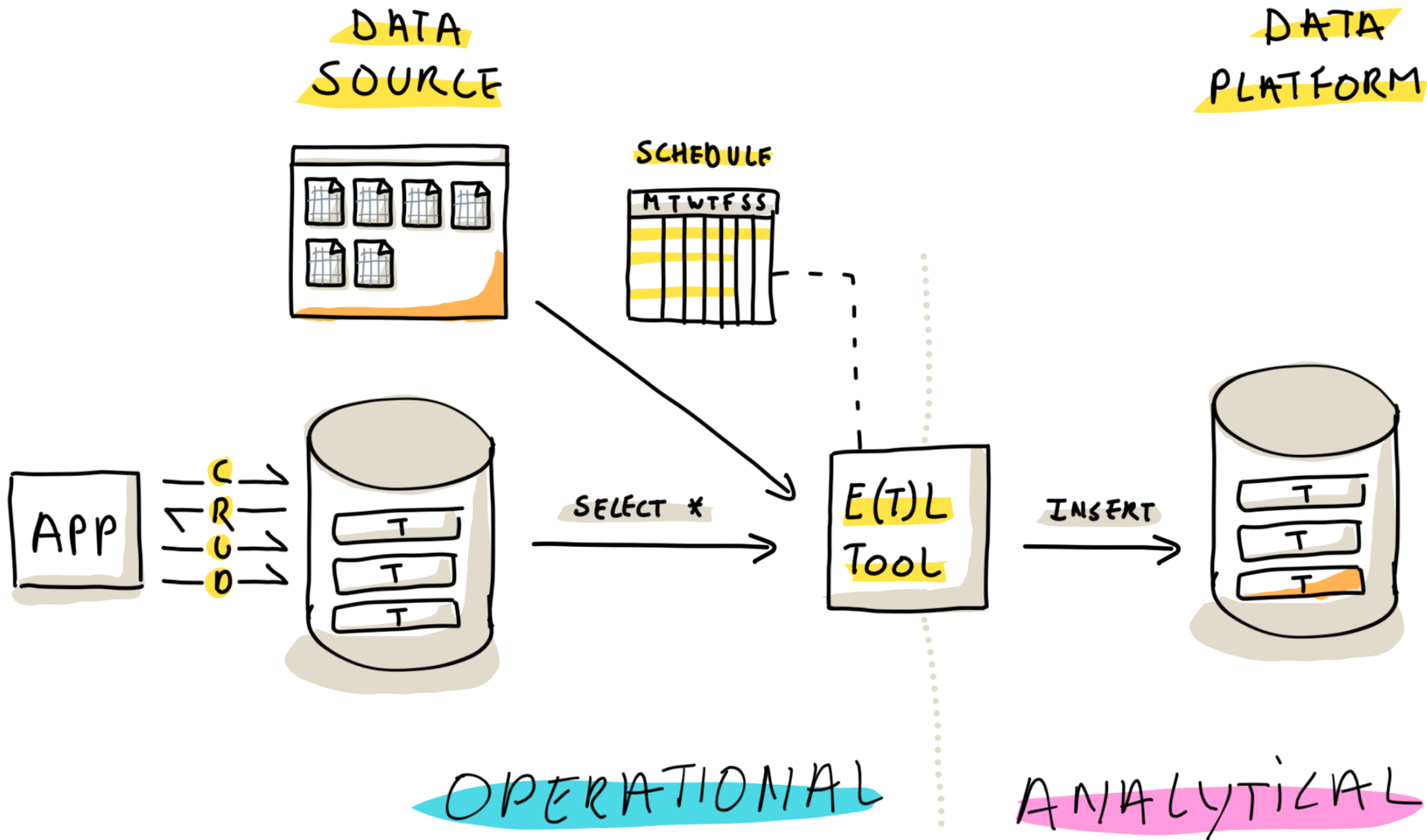
Flavor #4
Custom Builds



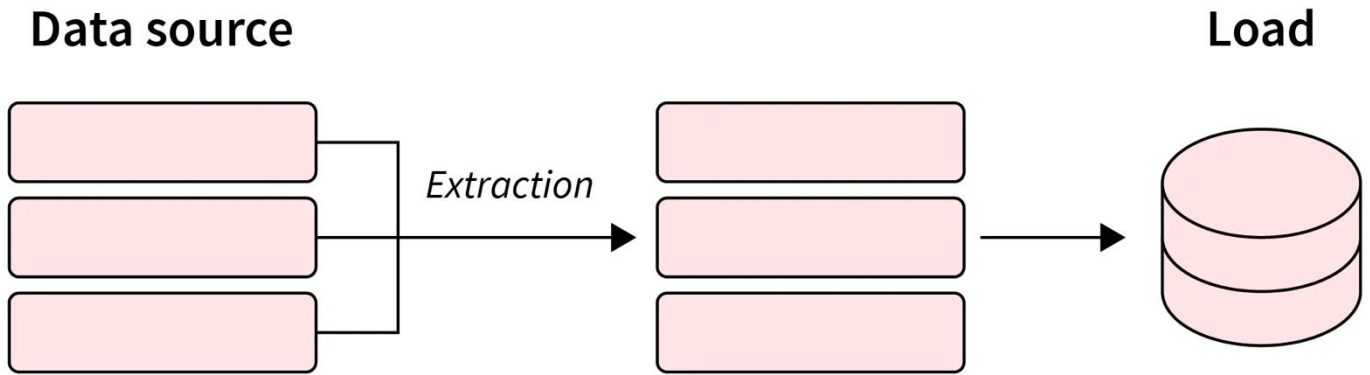
Flavor #1

Batch Loading

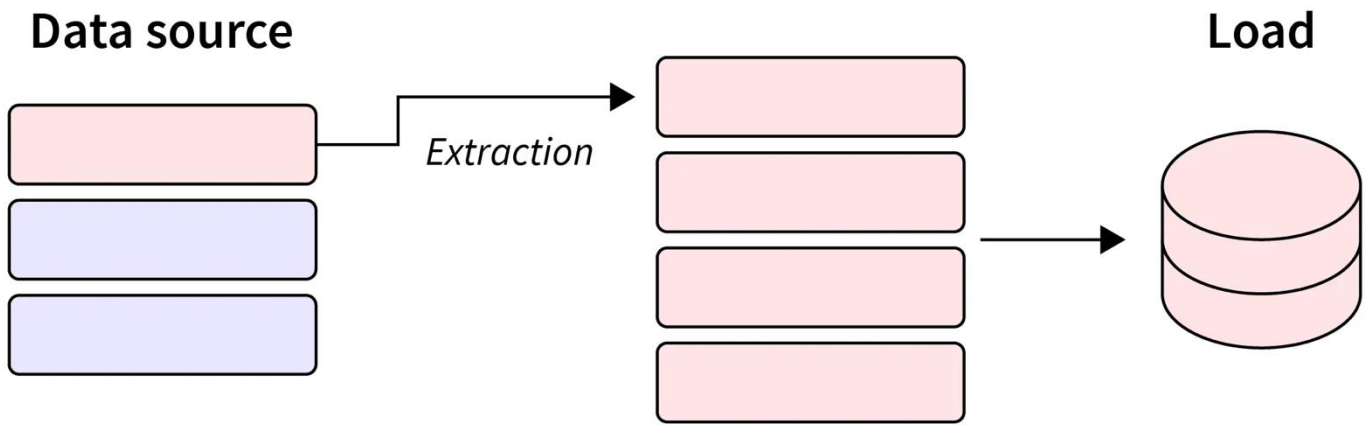




Full load
All available data is extracted at the same time























Incremental load
The occurring changes in the source data is incrementally extracted and loaded



EXERCISE: HOW CAN YOU INCREMENTALLY LOAD THIS TABLE?

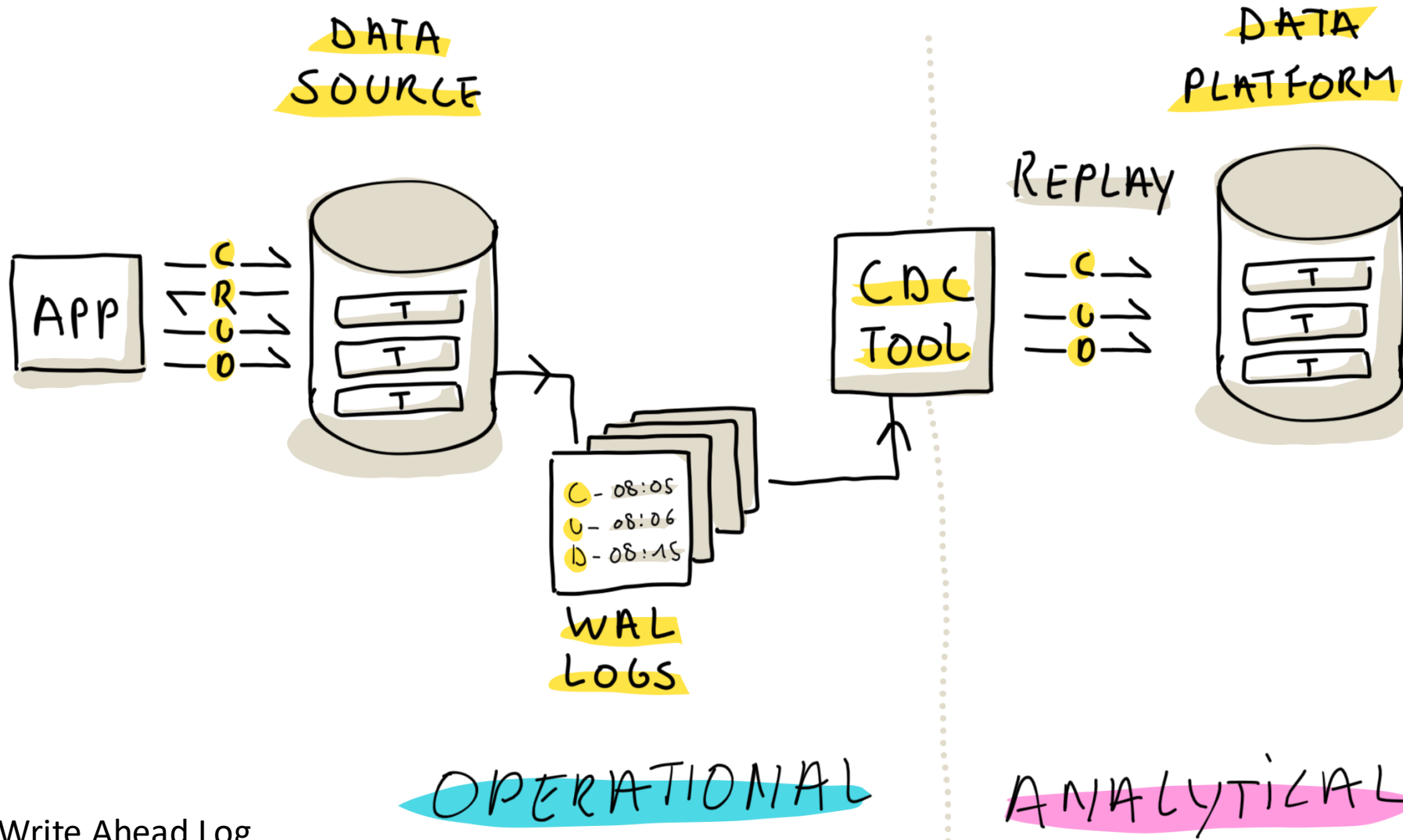
☒ m_transaction

OFF                    

	m_transaction_id	ad_client_id	ad_org_id	isactive	created	createdby	updated
1	1002856	1000010	1000061	Y	2020-03-18 14:40:01.253	100	2020-03-18 14:40:01.253
2	1002857	1000010	1000061	Y	2020-03-18 14:53:40.809	100	2020-03-18 14:53:40.809
3	1002858	1000010	1000061	Y	2020-03-18 15:20:00.275	100	2020-03-18 15:20:00.275
4	1002859	1000010	1000061	Y	2020-03-18 17:42:27.395	1000405	2020-03-18 17:42:27.395
5	1002860	1000010	1000129	Y	2020-03-18 19:50:49.07	100	2020-03-18 19:50:49.07
6	1002861	1000010	1000129	Y	2020-03-18 19:59:42.211	100	2020-03-18 19:59:42.211
7	1002862	1000010	1000129	Y	2020-03-18 20:00:39.243	100	2020-03-18 20:00:39.243
8	1002863	1000010	1000129	Y	2020-03-18 20:02:20.357	100	2020-03-18 20:02:20.357
9	1002864	1000010	1000129	Y	2020-03-18 20:10:12.598	100	2020-03-18 20:10:12.598
10	1002865	1000010	1000061	Y	2020-03-18 20:25:26.384	100	2020-03-18 20:25:26.384

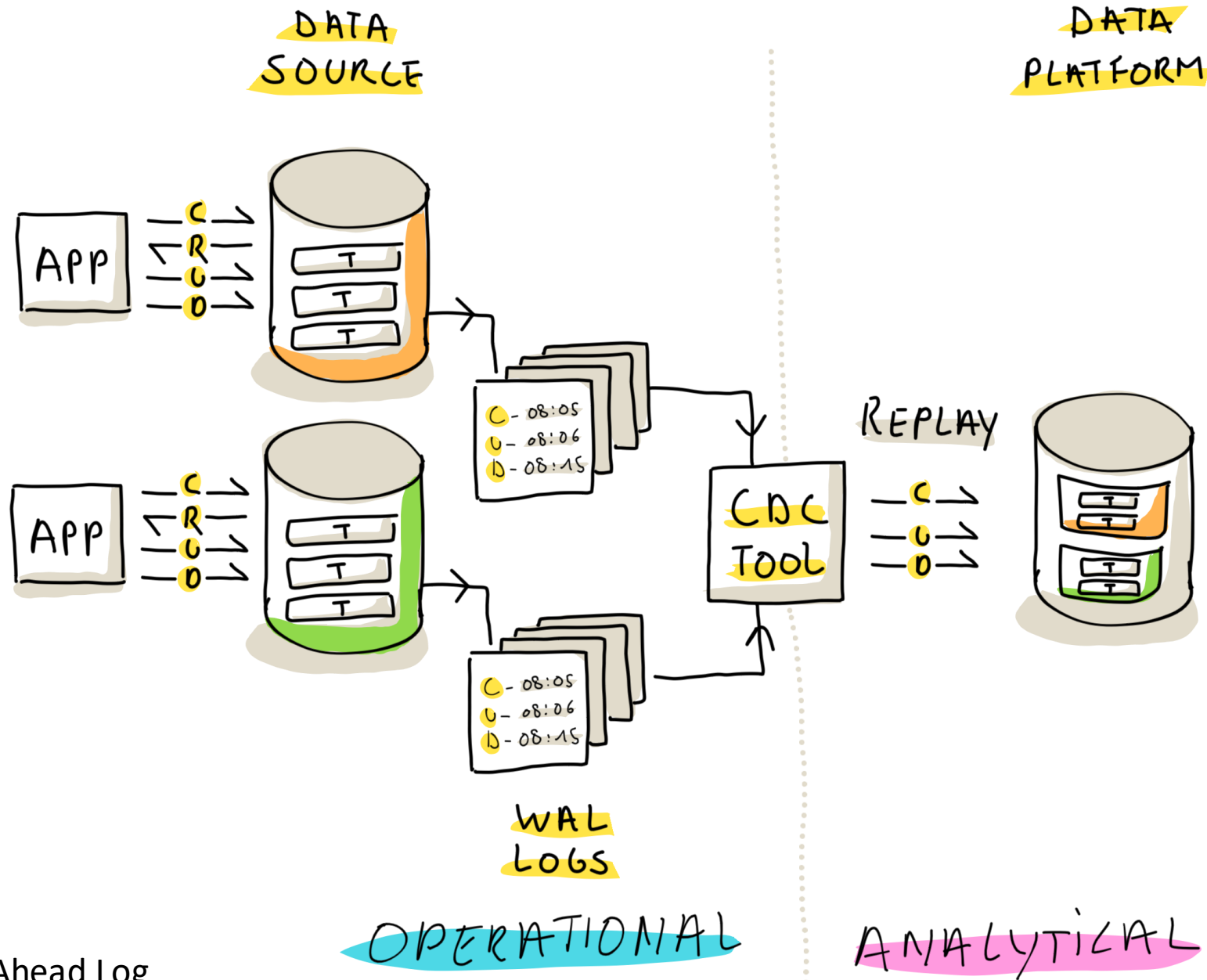
Flavor #2

CDC – Change Data Capture



WAL = Write Ahead Log

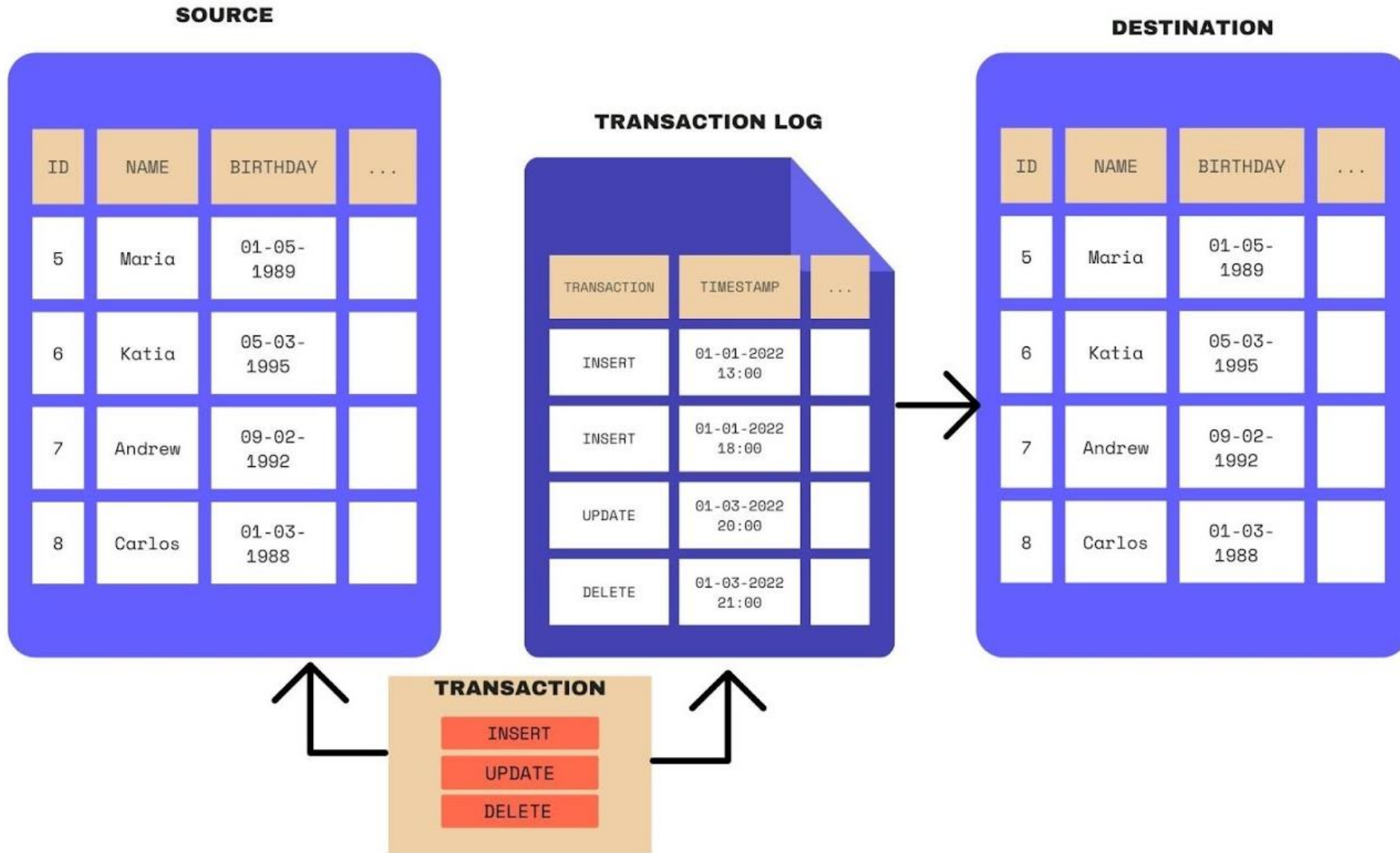


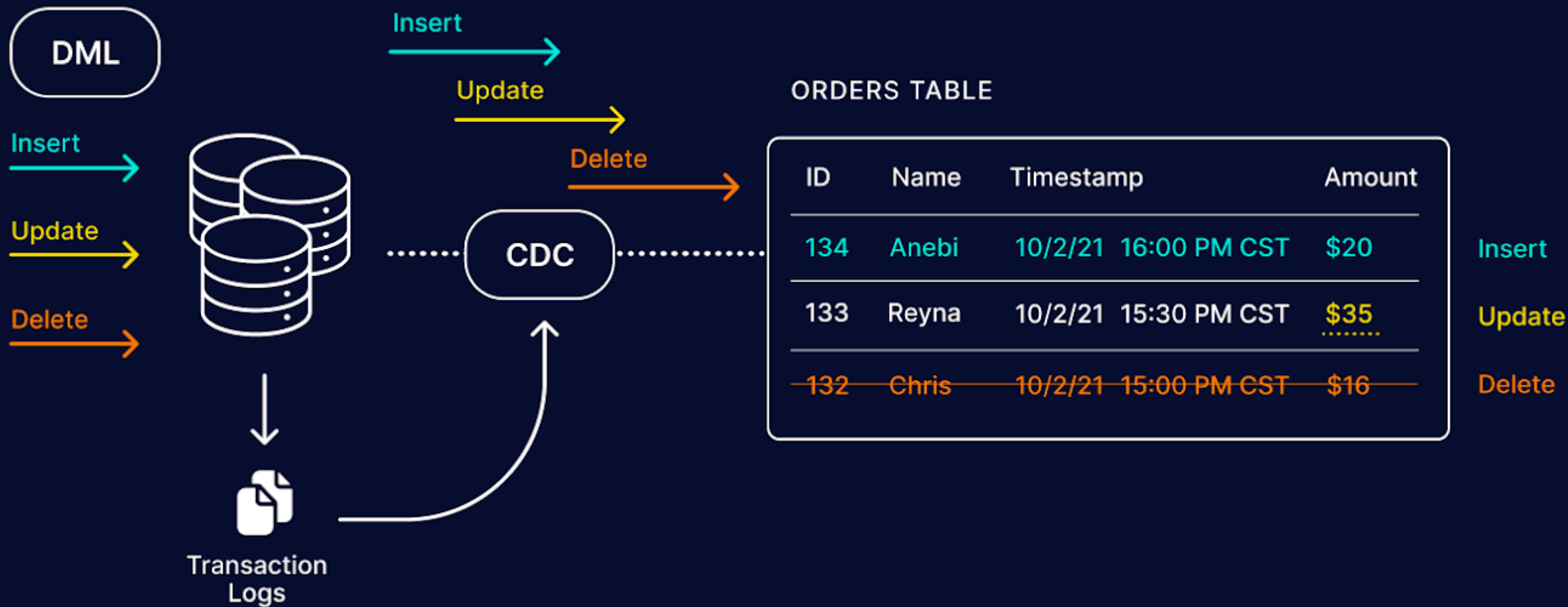


WAL = Write Ahead Log



Log-based CDC technique





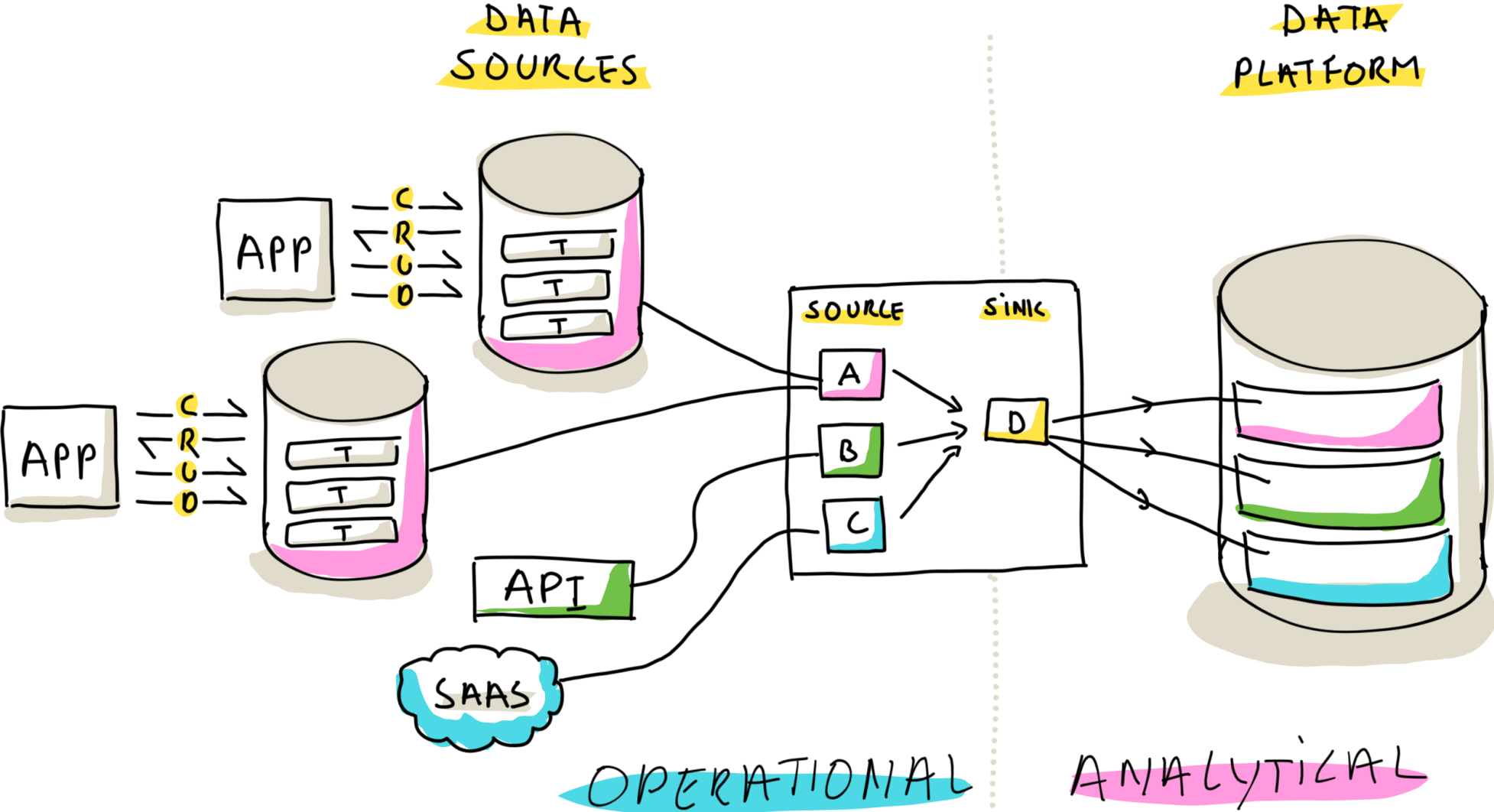
Change Data Capture (CDC) = a software that allows detecting and capturing changes made to data in a database and sending these changes, sometimes in real-time, to a downstream process or system. More specifically, CDC entails recording INSERT, UPDATE, and DELETE transactions applied to a table.

Various techniques exist: meta-data based, trigger based, **log based**

Log based CDC systems read data directly from the database Change Data Capture logs to identify changes in a database (not from the actual database)



Flavor #3 Connector Based



Connectors

34 Active - 12 Broken - 2 Paused • Last refreshed a day ago

ADD CONNECTOR

Connectors

48

Search by name...

All sources

All statuses

Transformations

67

Uploads

Destination

Logs

Users

Alerts

14

Notifications

Docs

Status

Name	Source	Status	Last synced
sql_server	SQL Server RDS	ACTIVE	a day ago
azure_function	Azure Functions	ACTIVE	a day ago
ss_demo	SQL Server RDS	ACTIVE	a day ago
salesforce_sandbox_sa...	Salesforce sandbox	ACTIVE	a day ago
gcs.customer	Google Cloud Sto...	ACTIVE	a day ago
gsheets.sales	Google Sheets	ACTIVE	a day ago
pg	Google Cloud Pos...	ACTIVE	a day ago
github	GitHub	ACTIVE	a day ago
netsuite	NetSuite SuiteAn...	ACTIVE	a day ago
fivetran_log	Fivetran Log	ACTIVE	a day ago
salesforce_sandbox_45...	Salesforce sandbox	ACTIVE	a day ago
adwords	Google Ads (AdW...	ACTIVE	a day ago
salesforce_sandbox	Salesforce sandbox	ACTIVE	a day ago
dark_sky	Google Cloud Fun...	ACTIVE	a day ago

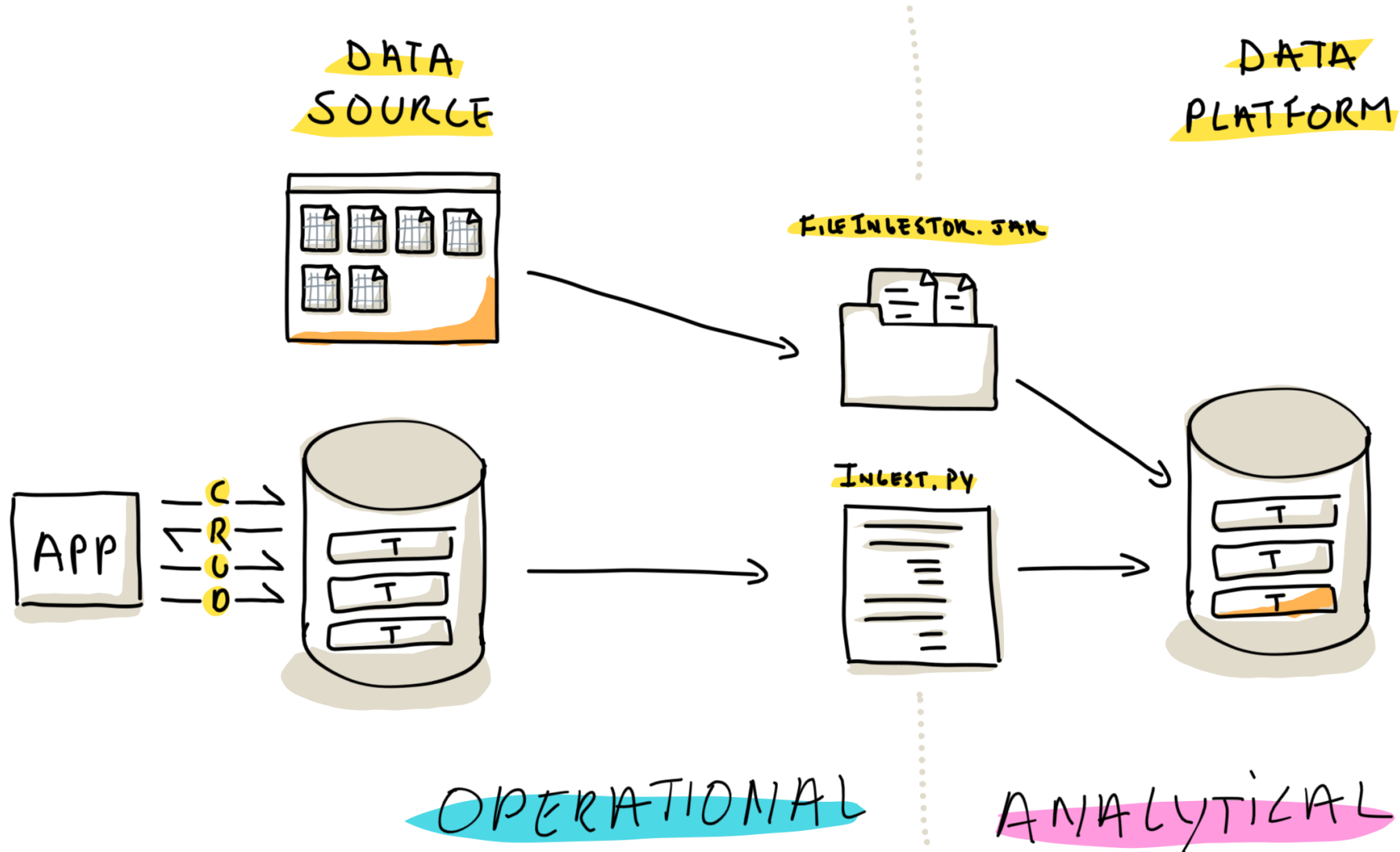


Pre-built source connectors are coupled with pre-built sink connectors, often using a graphical user interface.

- **Source connectors:** can be very diverse databases, APIs, SAAS-applications, applications, files, ...
- **Sink connectors:** mostly limited to (analytical) databases or data lakes.
- Highly **flexible** but **No control** over the individual connectors



Flavor #4 Custom Builds



Disadvantages:

- Building ingestion pipelines is usually more expensive than expected
- Specific programming knowledge & team needed
- High maintenance cost

Advantages:

- Full control
- Allows to ingest very specific / unique / exotic data sources



Building

VS

Buying

✓ **Customization and scale**

✓ **Greater control**
No license fees

✓ **Competitive edge**

✓ **Easy to modify**

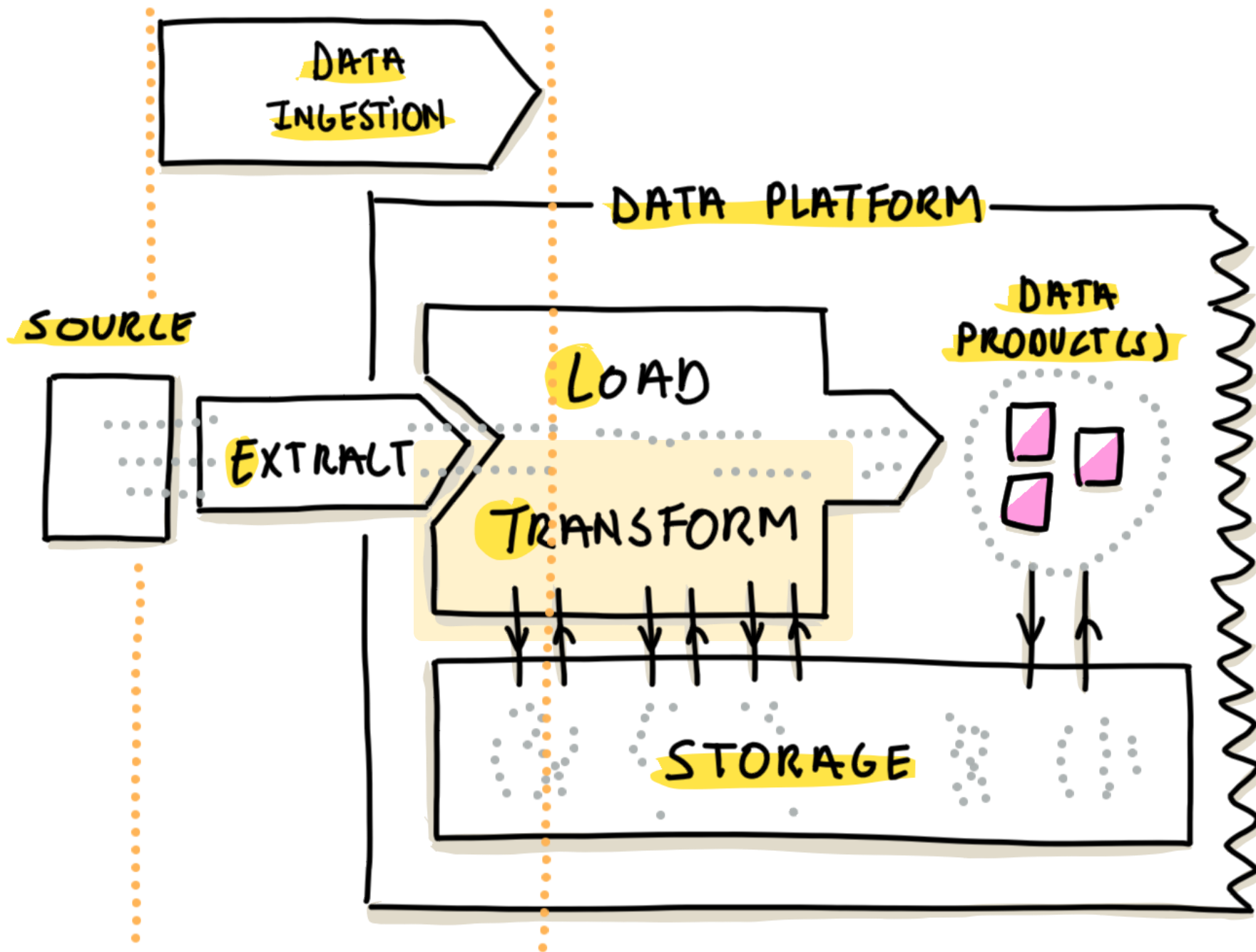
✓ **Lower upfront cost**

✓ **Rapid deployment**

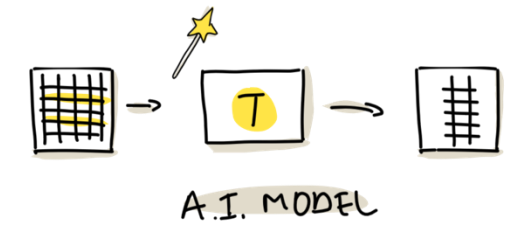
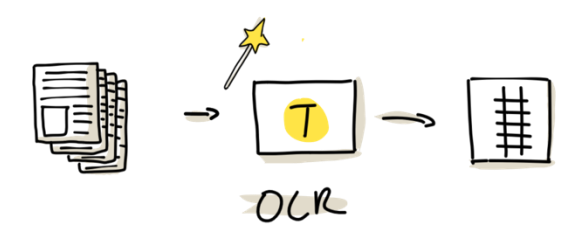
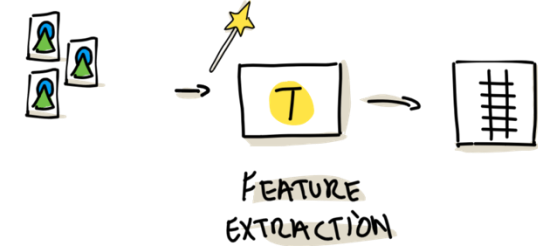
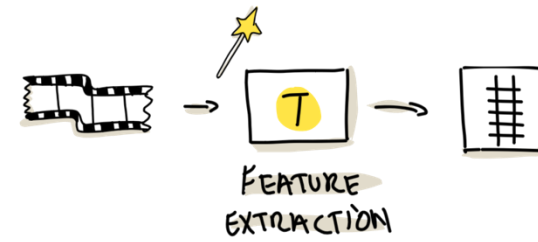
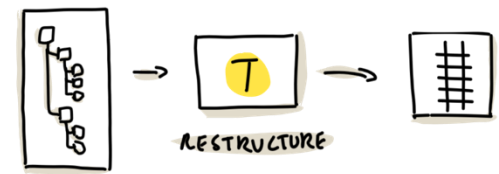
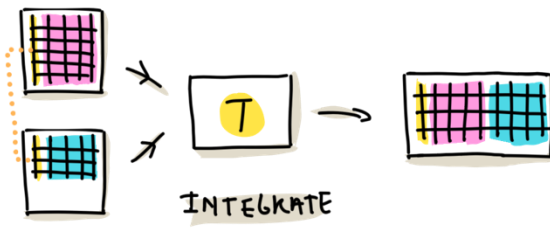
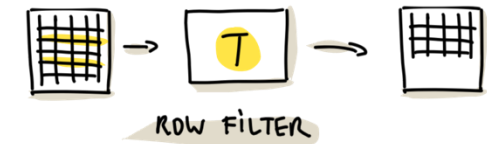
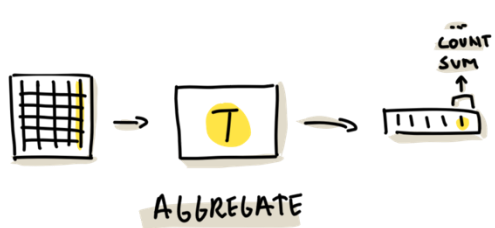
✓ **Updates and maintenance**

✓ **Has active userbase**

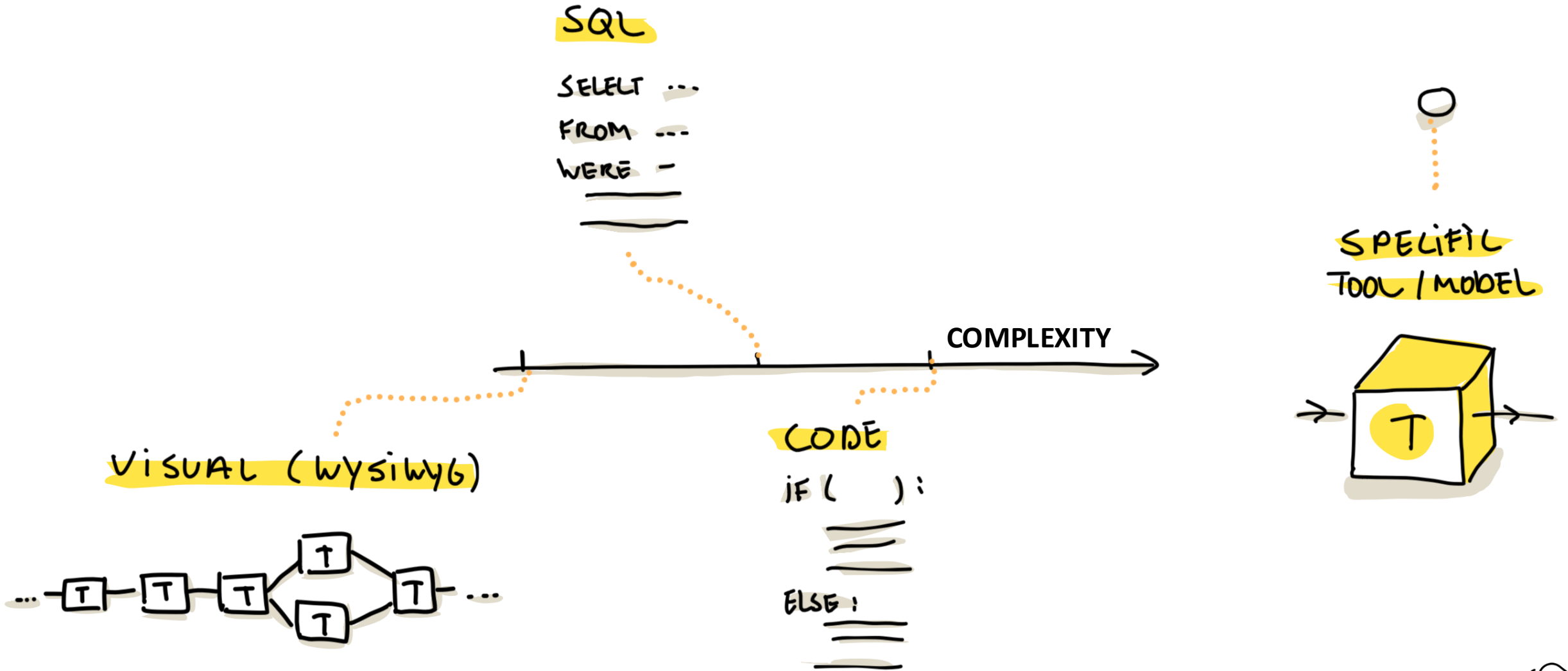




Data Transformations



Data Transformation Tool Flavors

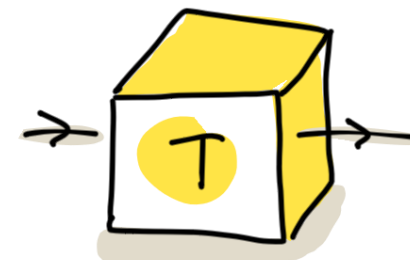


Data Transformation Tool Flavors

SQL

```
SELECT ...  
FROM ...  
WHERE -  
=====
```

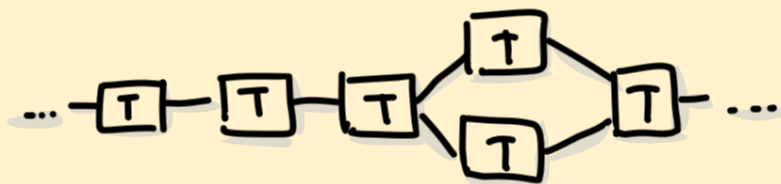
SPECIFIC
TOOL / MODEL



COMPLEXITY



VISUAL (WYSIWYG)

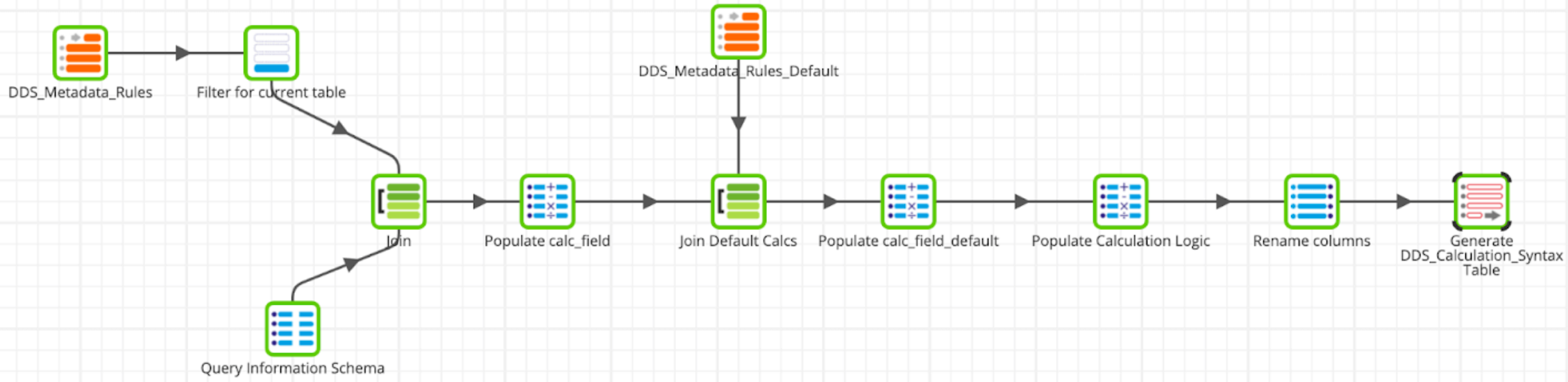


CODE

```
IF ( ):  
=====  
ELSE:  
=====  
=====
```

Purpose:

The purpose of this transformation step is to formulate the appropriate syntax that should be applied onto the columns during the transformation step of the process. All of the syntax will be dynamically populated into a syntax metadata table called 'DDS_Calculation_Syntax'.



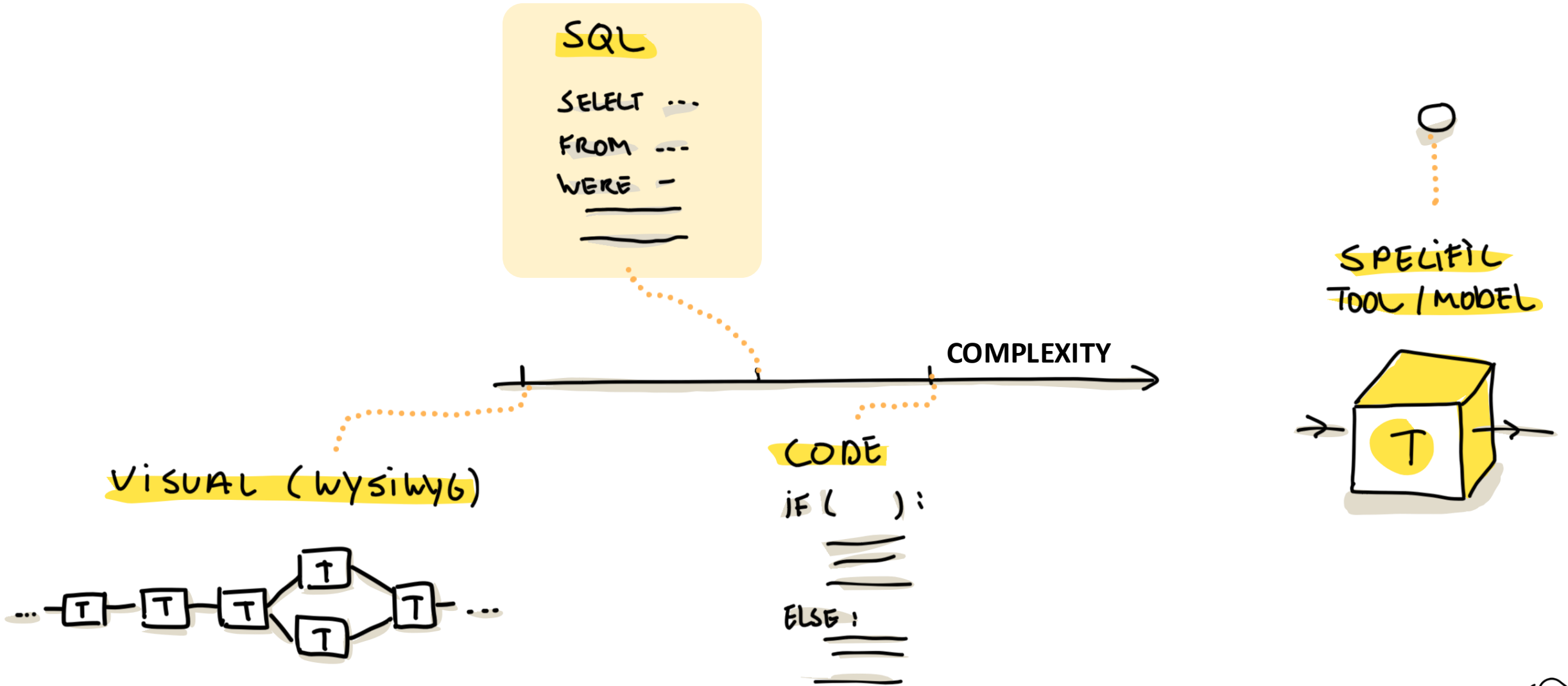
Properties | Export | Sample | Metadata | SQL | Plan | Help

Data | Row Count | Filter Not Set | Export

SCHEMA	TABLE_NAME	COLUMN_NAME	calculation
KB_AZURE_SNOW	DDS_STATE	STATEID	STATEID
KB_AZURE_SNOW	DDS_STATE	COUNTRYNAME	INITCAP(NVL(TRIM("COUNTRYNAME"),' '))
KB_AZURE_SNOW	DDS_STATE	STATECODE	UPPER(NVL(TRIM("STATECODE"),' '))
KB_AZURE_SNOW	DDS_STATE	FLAG	TO_BOOLEAN("FLAG")
KB_AZURE_SNOW	DDS_STATE	DATE	TO_TIMESTAMP_NTZ(TO_VARCHAR("DATE"), 'yyyymmdd')
KB_AZURE_SNOW	DDS_STATE	STATENAME	INITCAP(NVL(TRIM("STATENAME"),' '))



Data Transformation Tool Flavors



Project

view docs ?

Scratchpad 1

fact_employee_detail.sql

open pull request...

branch: jbarcheski_dev_demo

dbt_generic_demo

analysis

data

dbt_modules

logs

macros

models

sources

staging

warehouse

human_resources

dim_department.sql

dim_department.yml

dim_employee_department.sql

dim_employee_department.yml

fact_employee_detail.sql

fact_employee_detail.yml

purchasing

snapshots

target

tests

.gitignore

dbt_project.yml

packages.yml

```

31
32 final as (
33
34     select
35         to_binary(hex_encode('businessentityid'), 'HEX') as employee_sk,
36         e.businessentityid,
37         e.nationalidnumber as national_id,
38         e.loginid as login_id,
39         e.organizationnode as organization_node,
40         e.organizationlevel as organization_level,
41         e.jobtitle as job_title,
42         e.birthdate as birth_date,
43         e.maritalstatus as martial_status,
44         e.gender as gender,
45         e.hiredate as hire_date,
46         e.salariedflag as salaried_flag,
47         e.vacationhours as vacation_hours,
48         e.sickleavehours as sick_leave_hours,
49         e.currentflag as employee_current_flag,
50         e.rowguid as row_guid,
51         e.modifieddate as employee_modified_date,
52         datediff(year, hiredate, current_date()) as years_since_hire,
53         current_department_id,
54         current_shift_id,
55         current_department_start_date
56     from employees e
57     left join current_department d
58         on e.businessentityid=d.businessentityid
59 )
60
61 select * from final

```

Preview

Compile

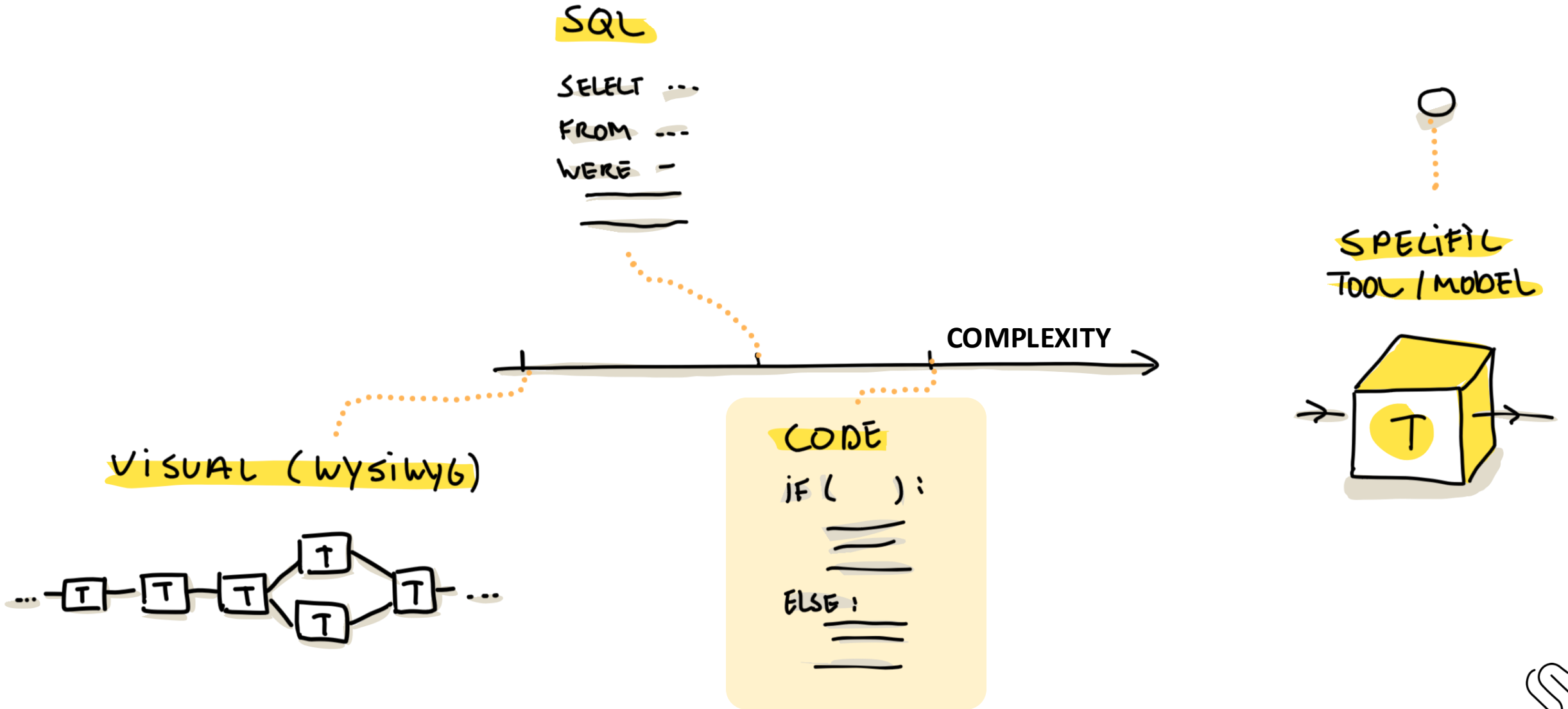
Query Results

Compiled SQL

Lineage



Data Transformation Tool Flavors



```
import time
from sklearn.datasets import make_classification
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression

# Set training and validation sets
X, y = make_classification(n_samples=1000000, n_features=1000, n_classes = 2)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=10000)

# Solvers
solvers = ['liblinear', 'saga']

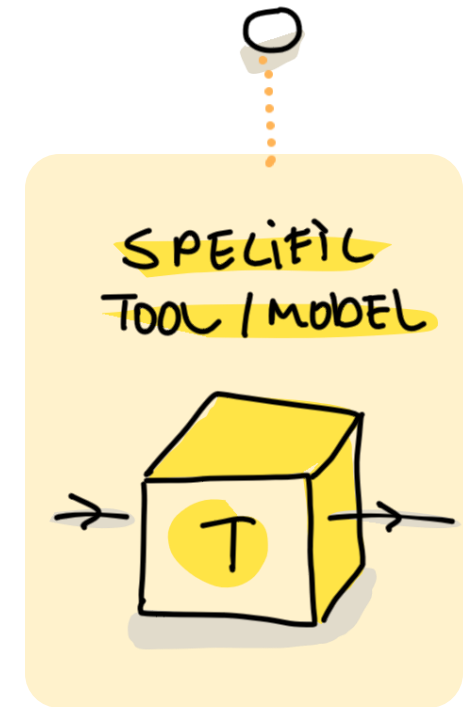
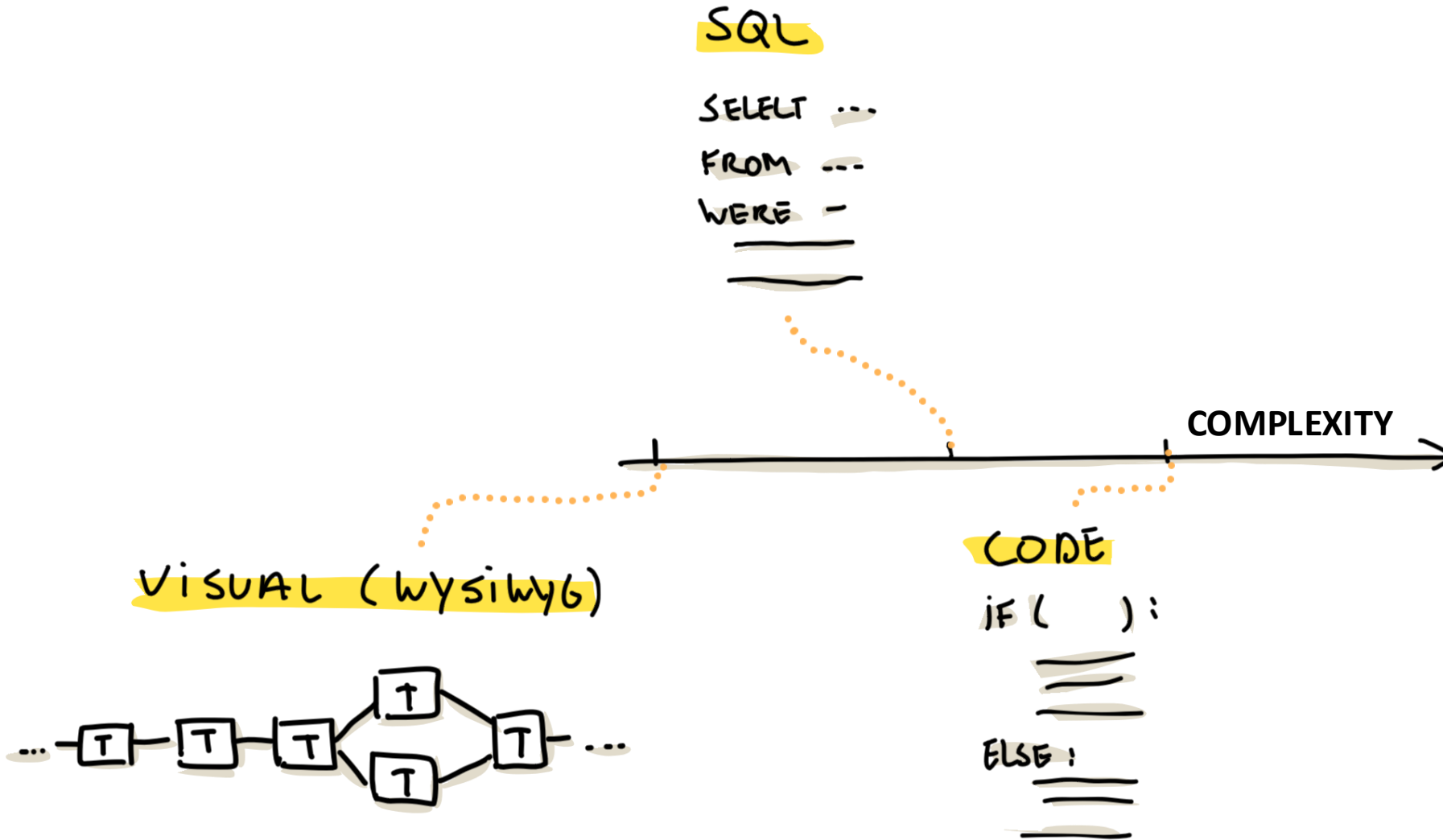
for sol in solvers:
    start = time.time()
    logreg = LogisticRegression(solver=sol)
    logreg.fit(X_train, y_train)
    end = time.time()
    print(sol + " Fit Time: ", end-start)
```

liblinear Fit Time: 2424.526051044464

saga Fit Time: 133.299968957901



Data Transformation Tool Flavors





Background-color change



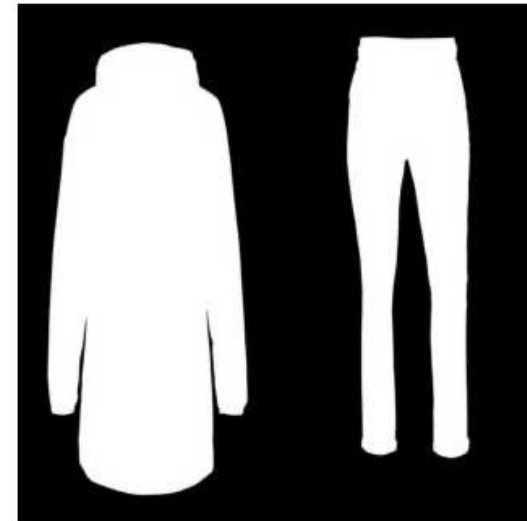
Background-color change

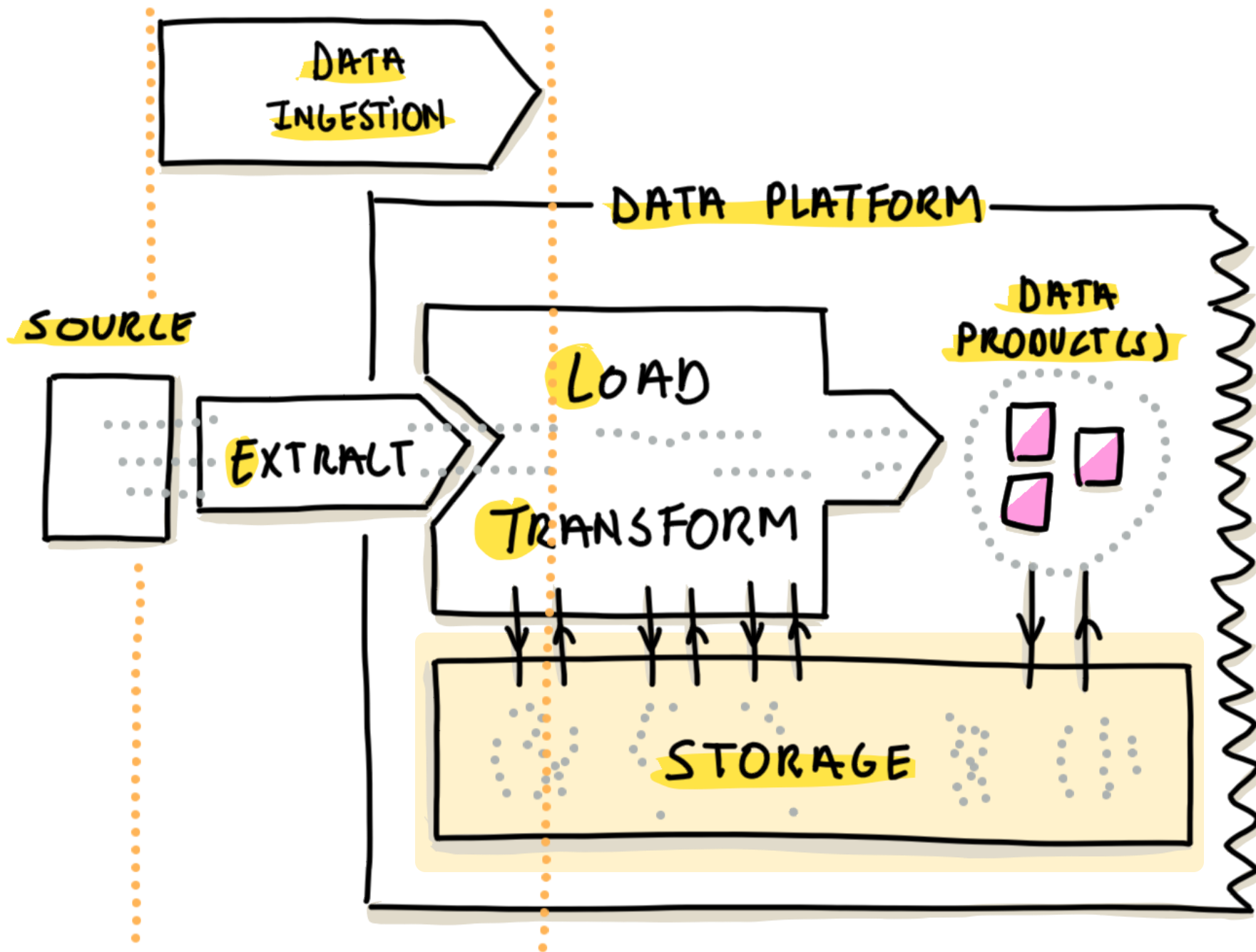


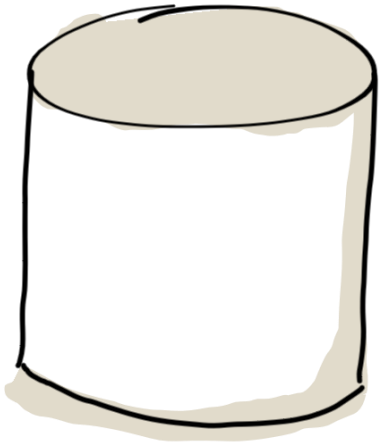
Cropping



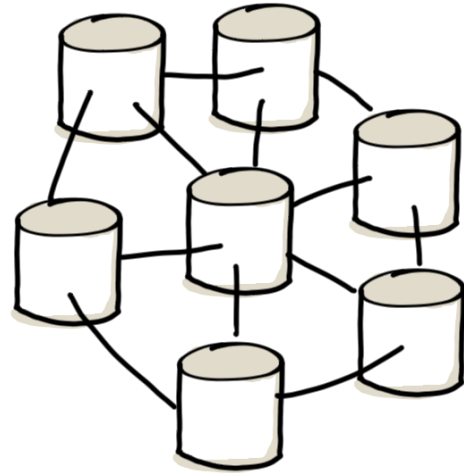
Combined images



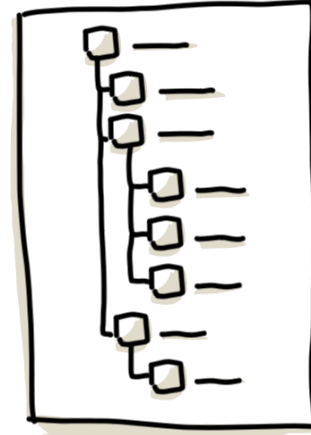




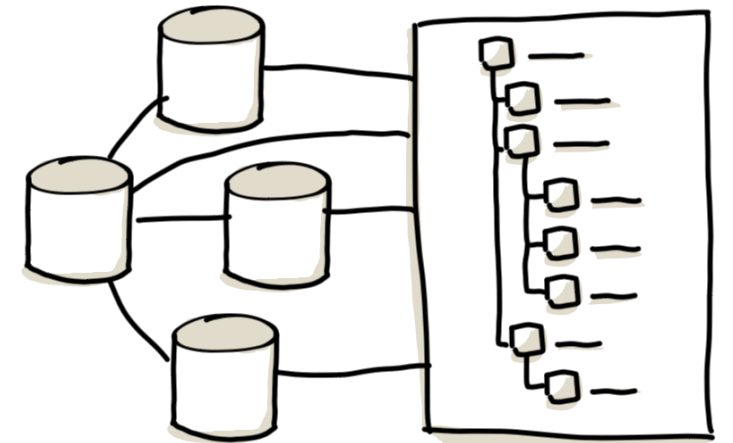
DATABASE



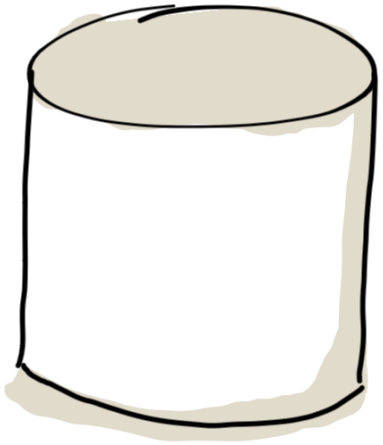
ANALYTICAL
DATABASE



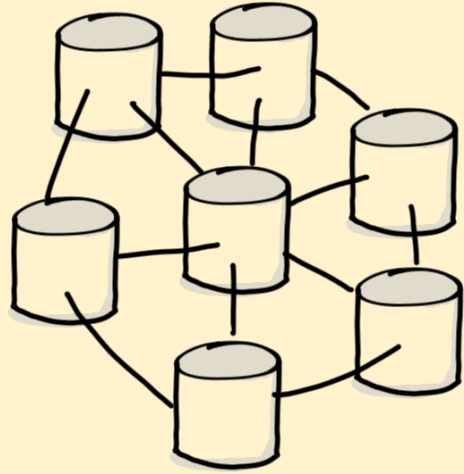
DATA LAKE



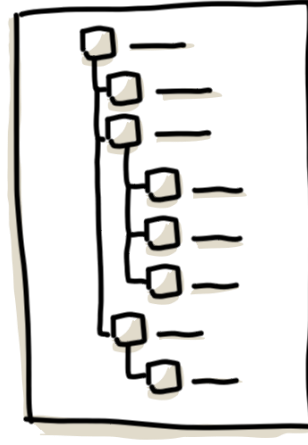
DATA
LAKEHOUSE



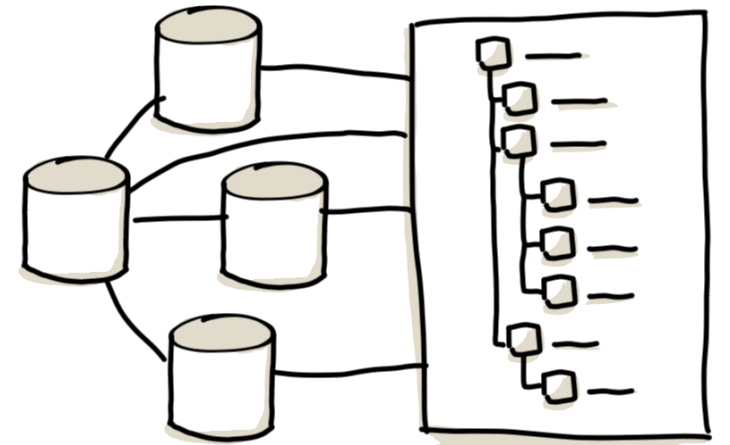
DATABASE



ANALYTICAL
DATABASE

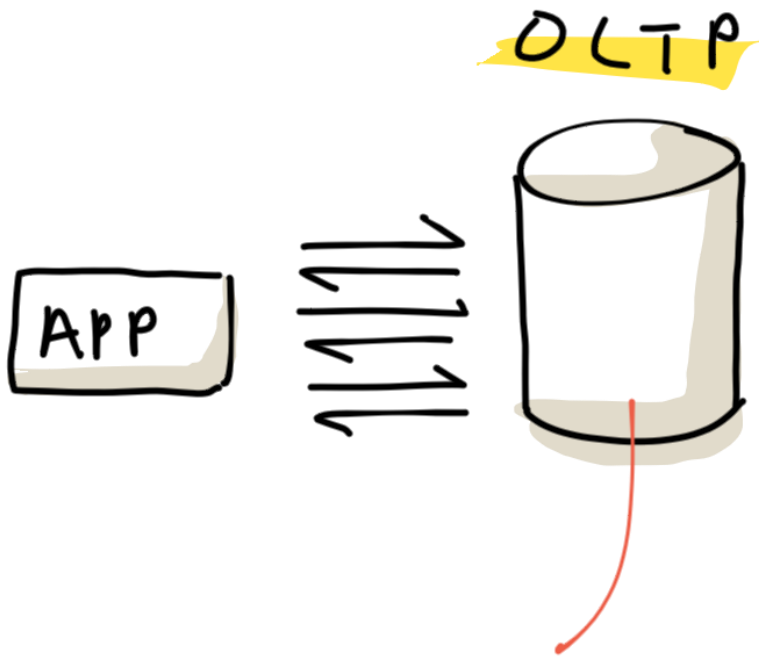


DATA LAKE



DATA
LAKEHOUSE

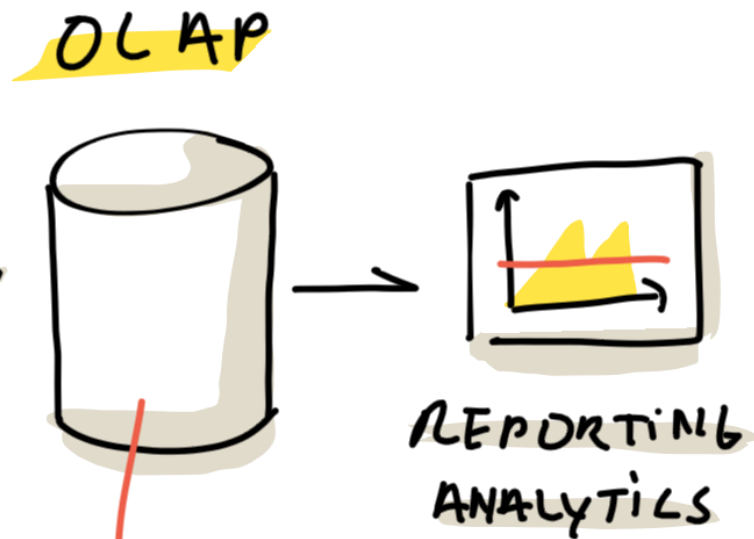
DATA SOURCES



TRANSACTIONAL
DATABASE

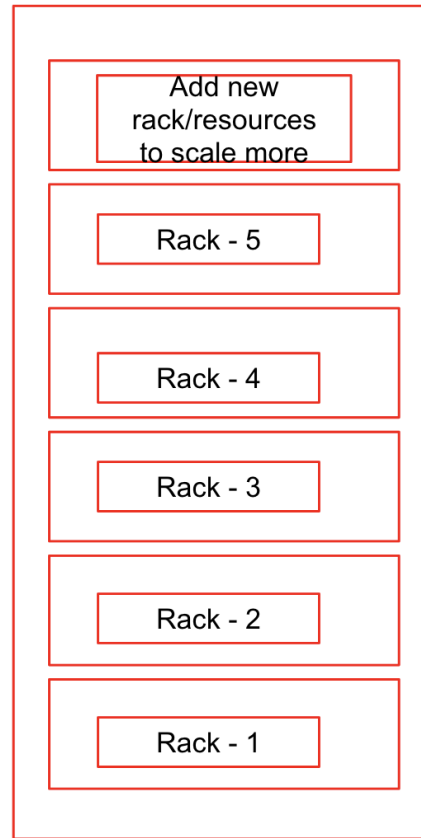


DATA PLATFORM



ANALYTICAL
DATABASE

A DWH Database (often called 'Cloud DWH') is tuned for **horizontal scaling**.



Host 1
192.168.1.1

Vertical Scaling

To scale more, Add more RAM, CPU, Memory to the **one existing machine**

Horizontal Scaling

To scale more: Add more machines to existing **group of distributed system**

Host 1
192.168.1.1

Host 2
192.168.1.2

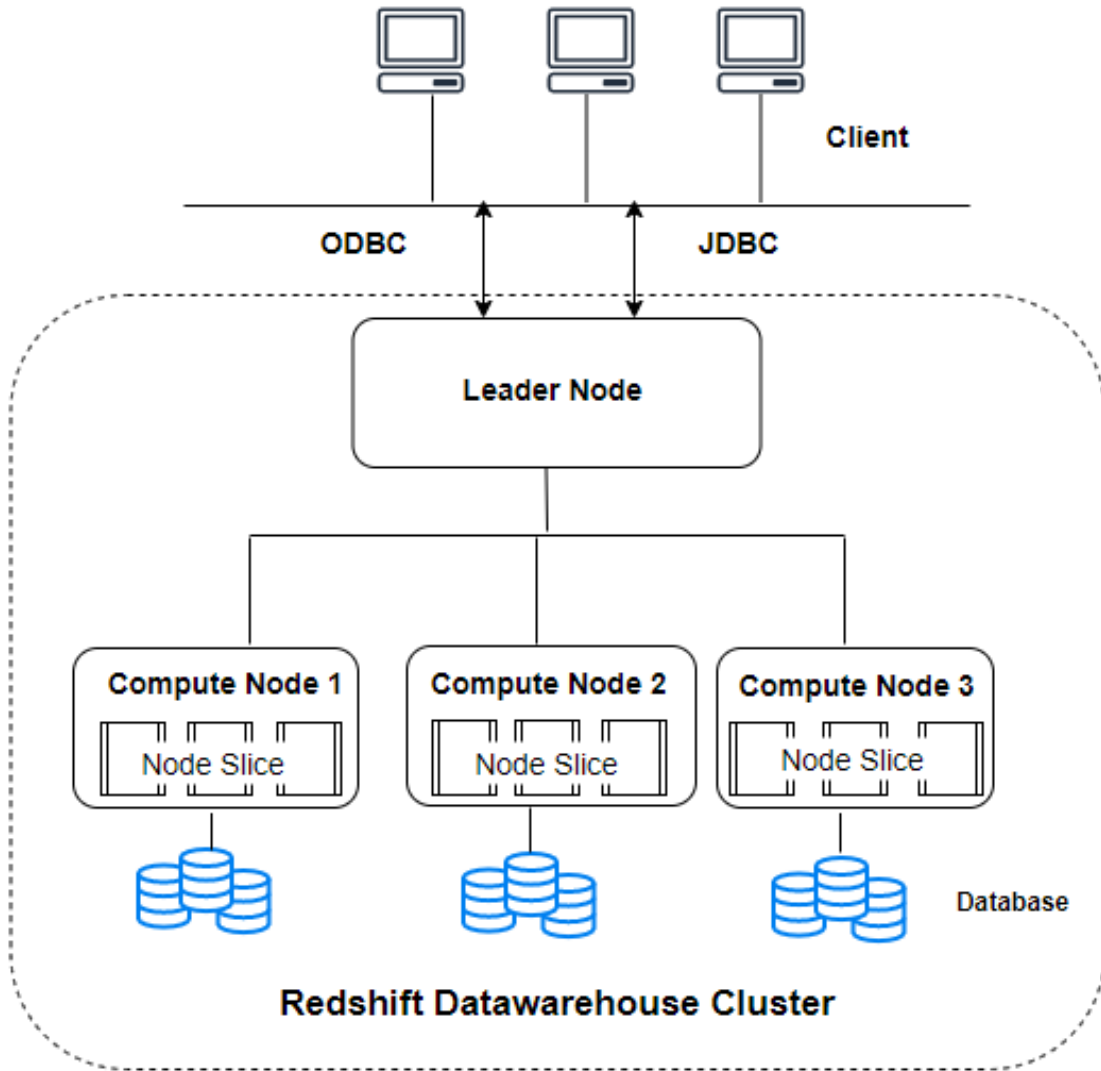
Host 3
192.168.1.3

Host x
192.168.1.x

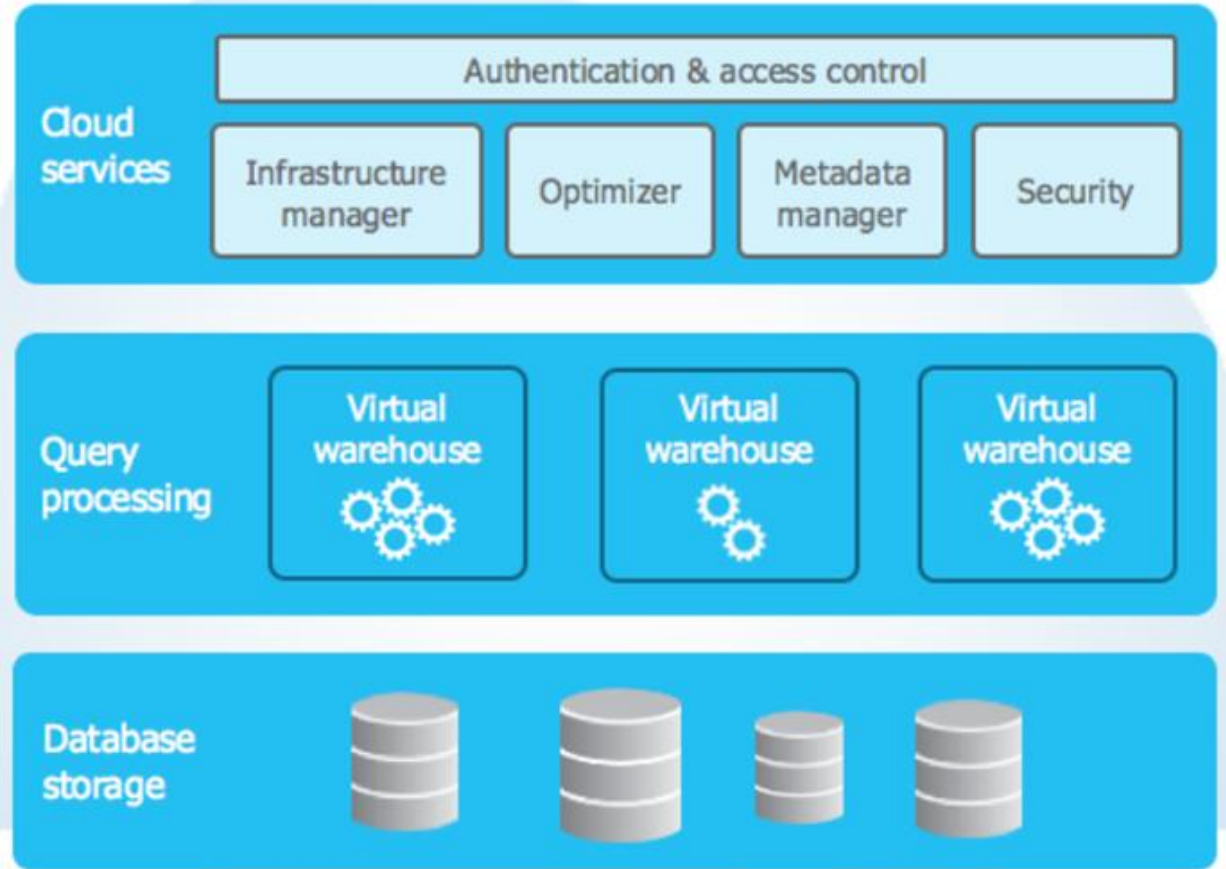
Add x+1
host to
scale out



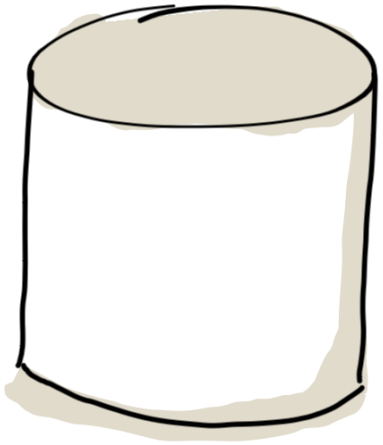
Example: Amazon Redshift



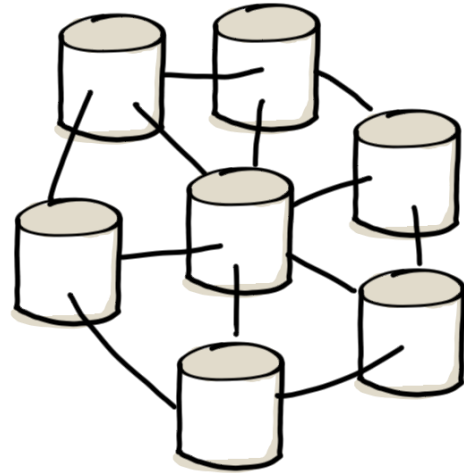
Example: Snowflake



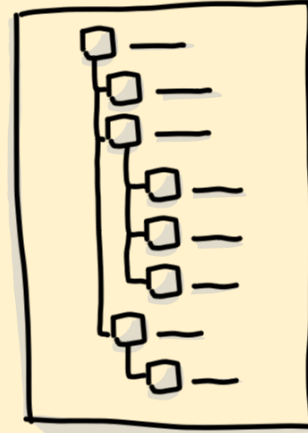




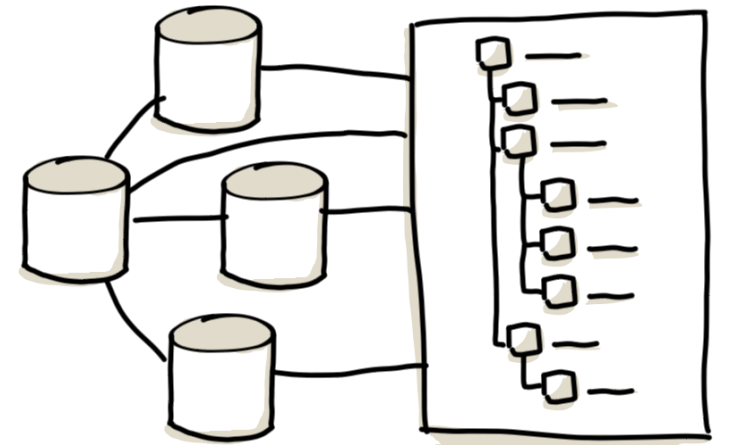
DATABASE



ANALYTICAL
DATABASE



DATA LAKE



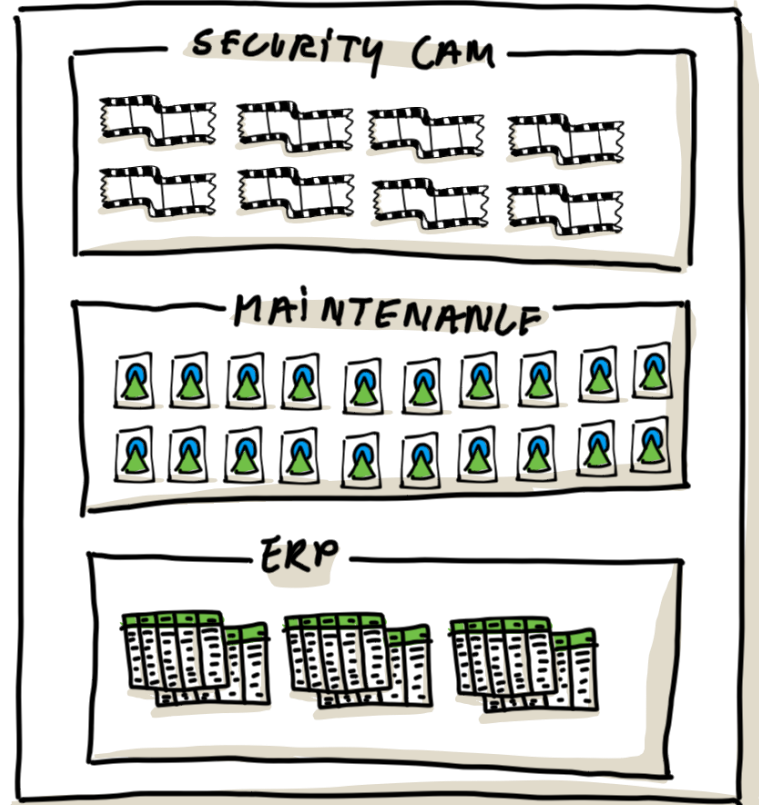
DATA
LAKEHOUSE

Data Lake

DATA SOURCES



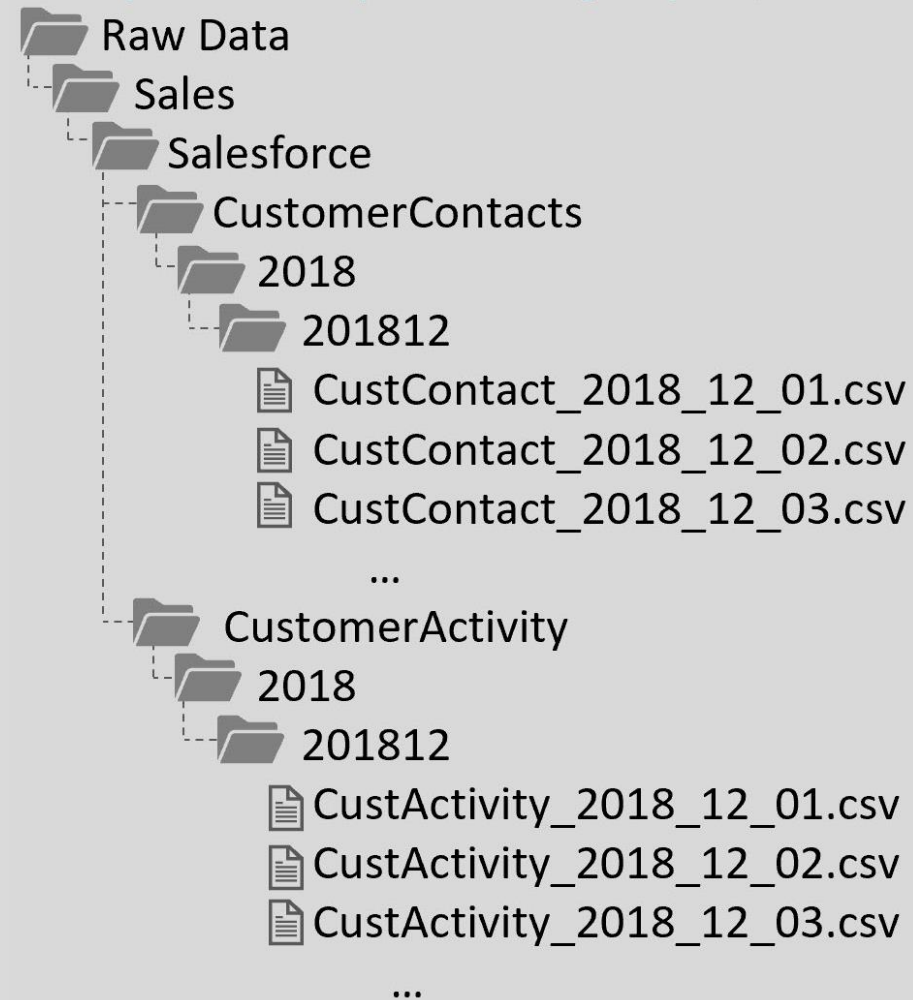
DATA LAKE

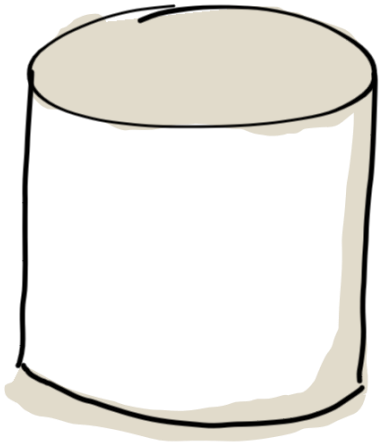


DATA PRODUCTS

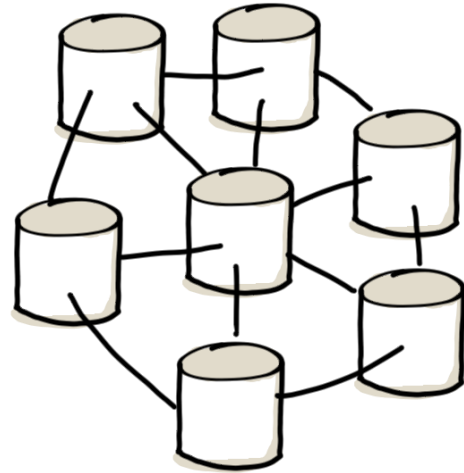


Data Lake: A **Structured** File Repository

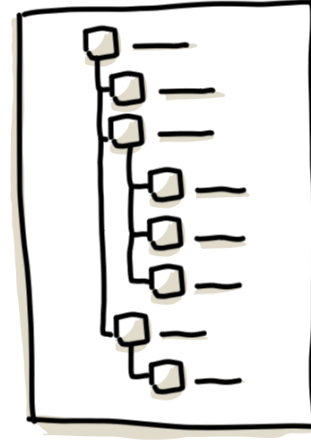




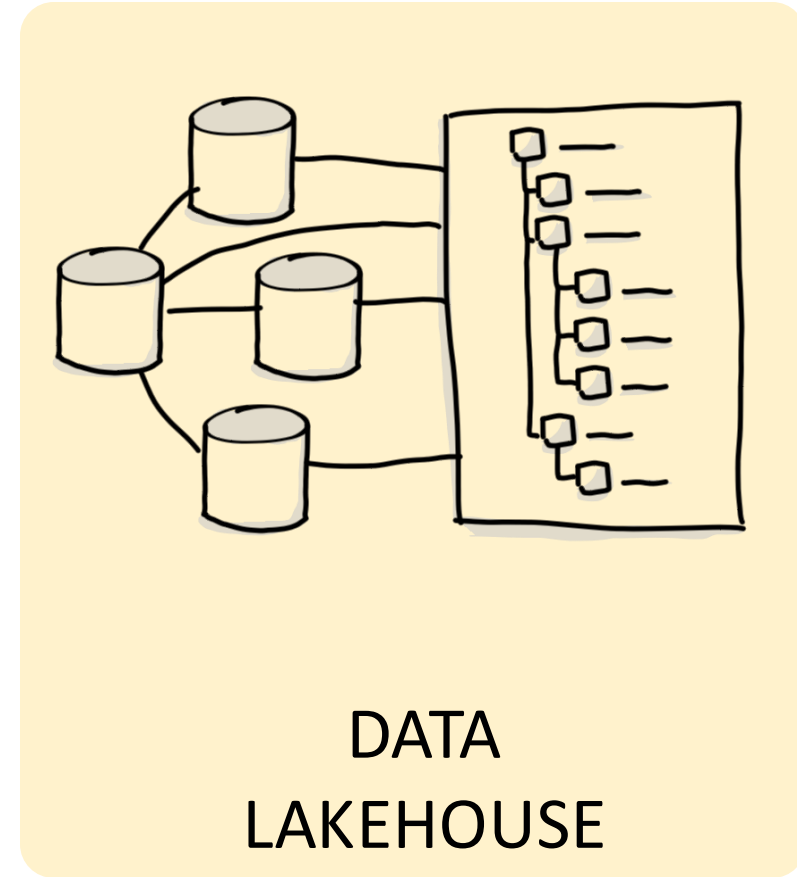
DATABASE



ANALYTICAL
DATABASE

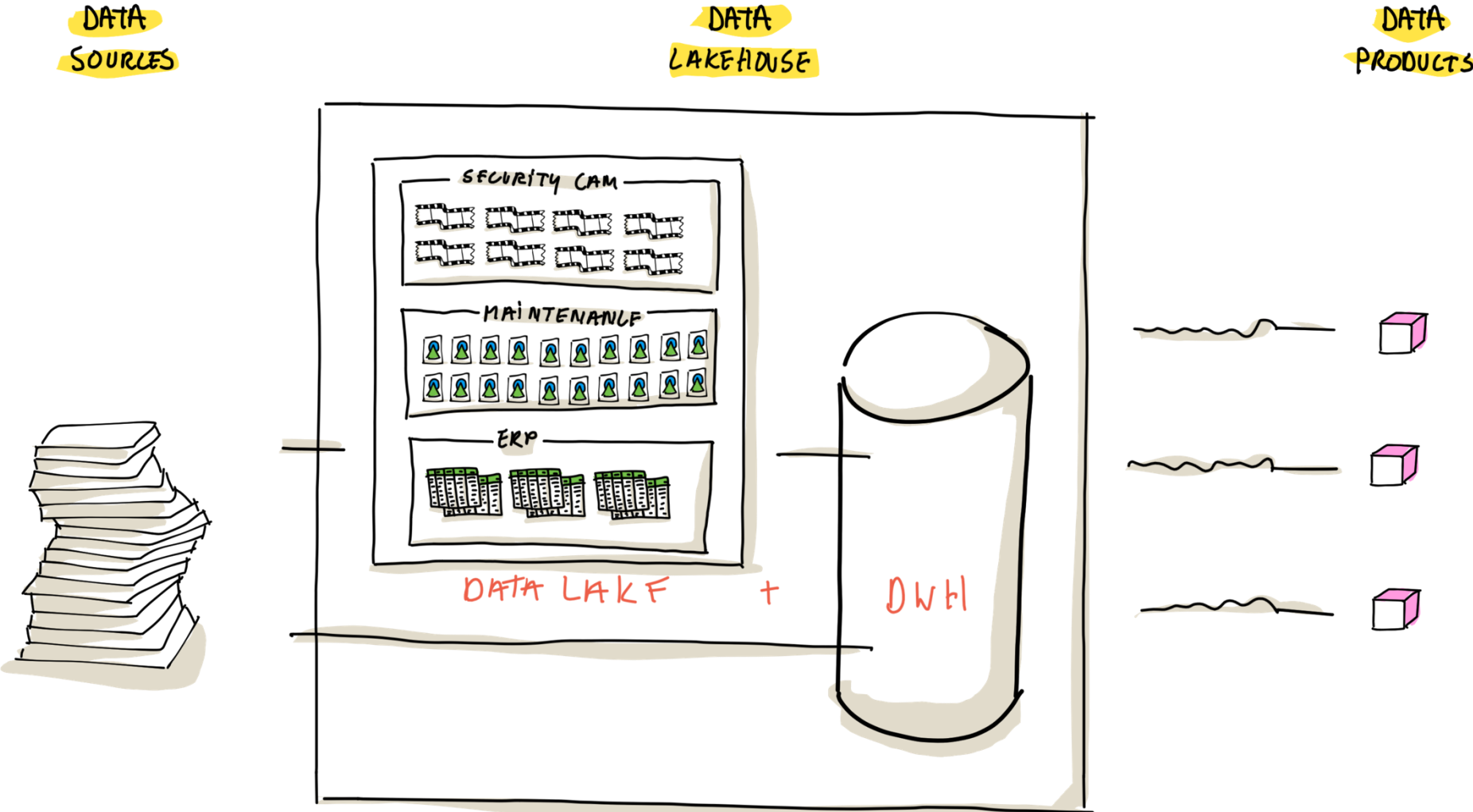


DATA LAKE



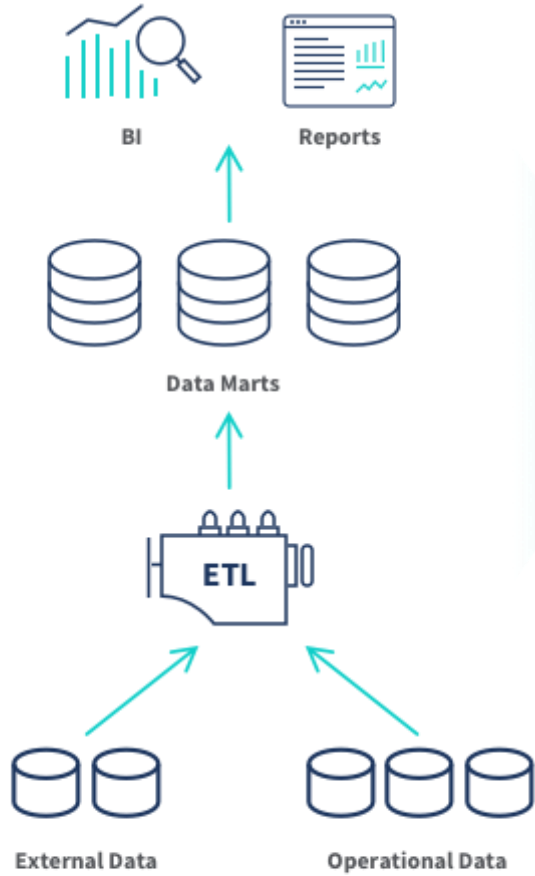
DATA
LAKEHOUSE

Data Lakehouse



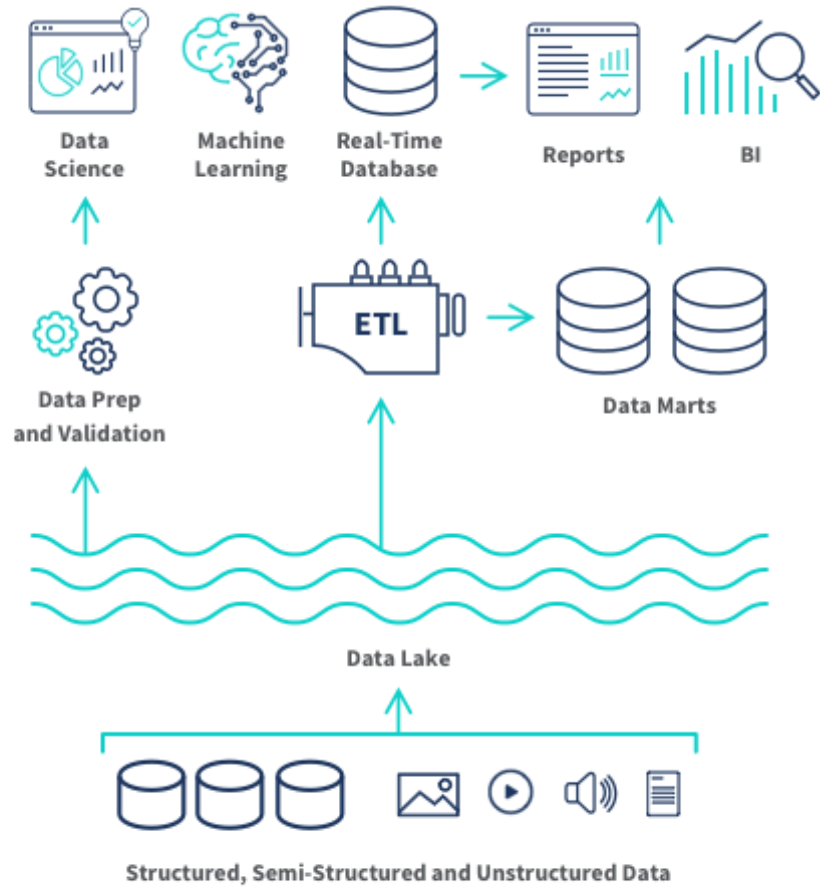
LATE 1980'S

Data Warehouse



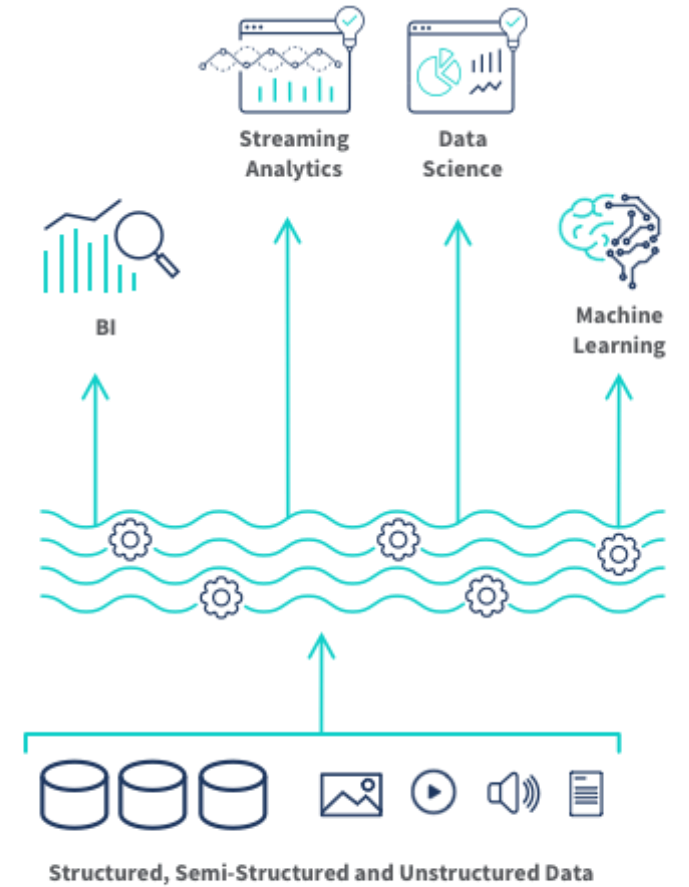
2011

Data Lake



2020

Lakehouse

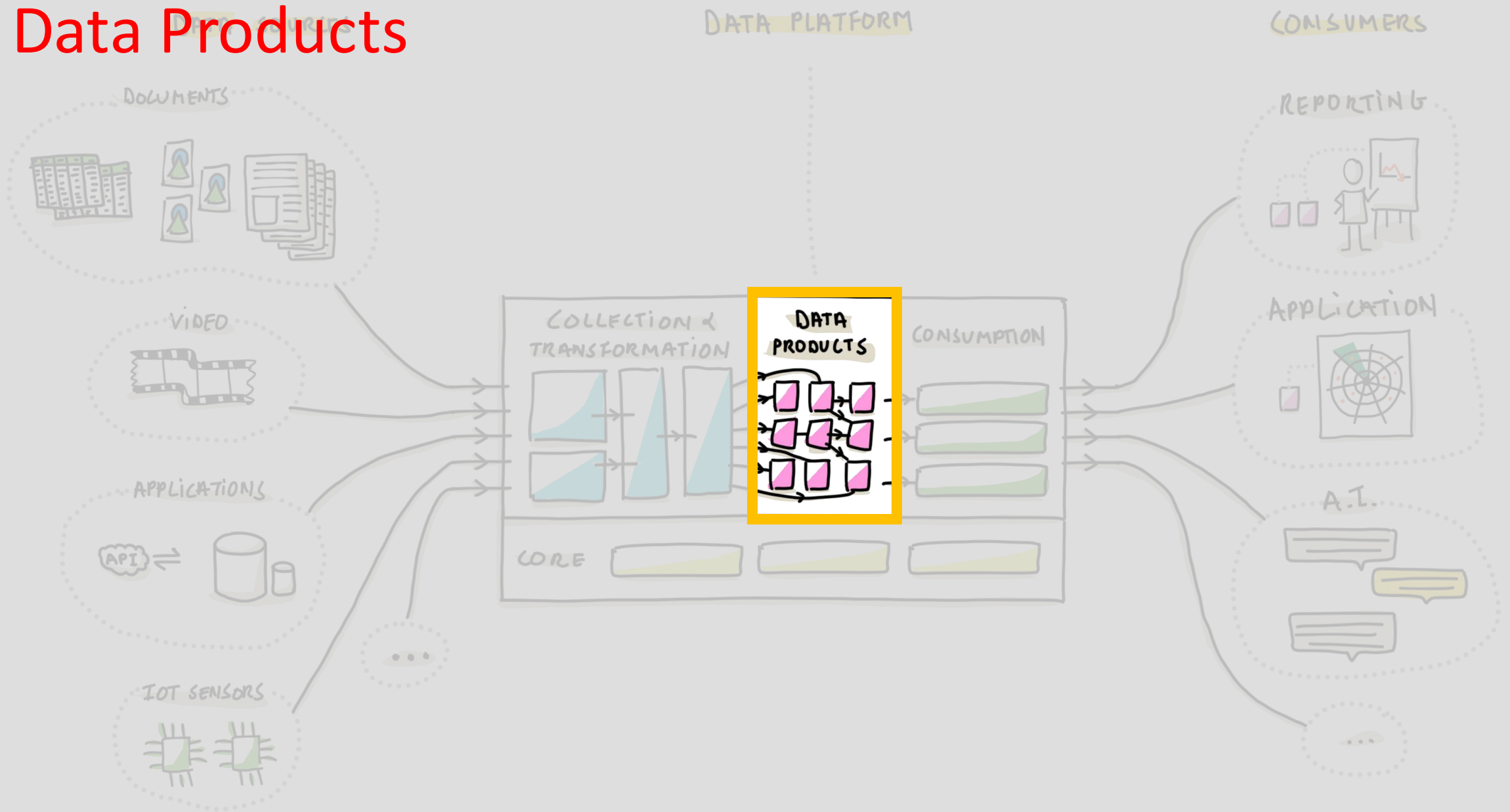


EXERCISE : DIFFERENCES BETWEEN STORAGE TYPES?

	Database	DWH Database	Data Lake	Data Lakehouse	?
Cost					
Scaling					
Volume					
Type of Data					
Performance					
Agility					
Users					
?					



2. Data Products



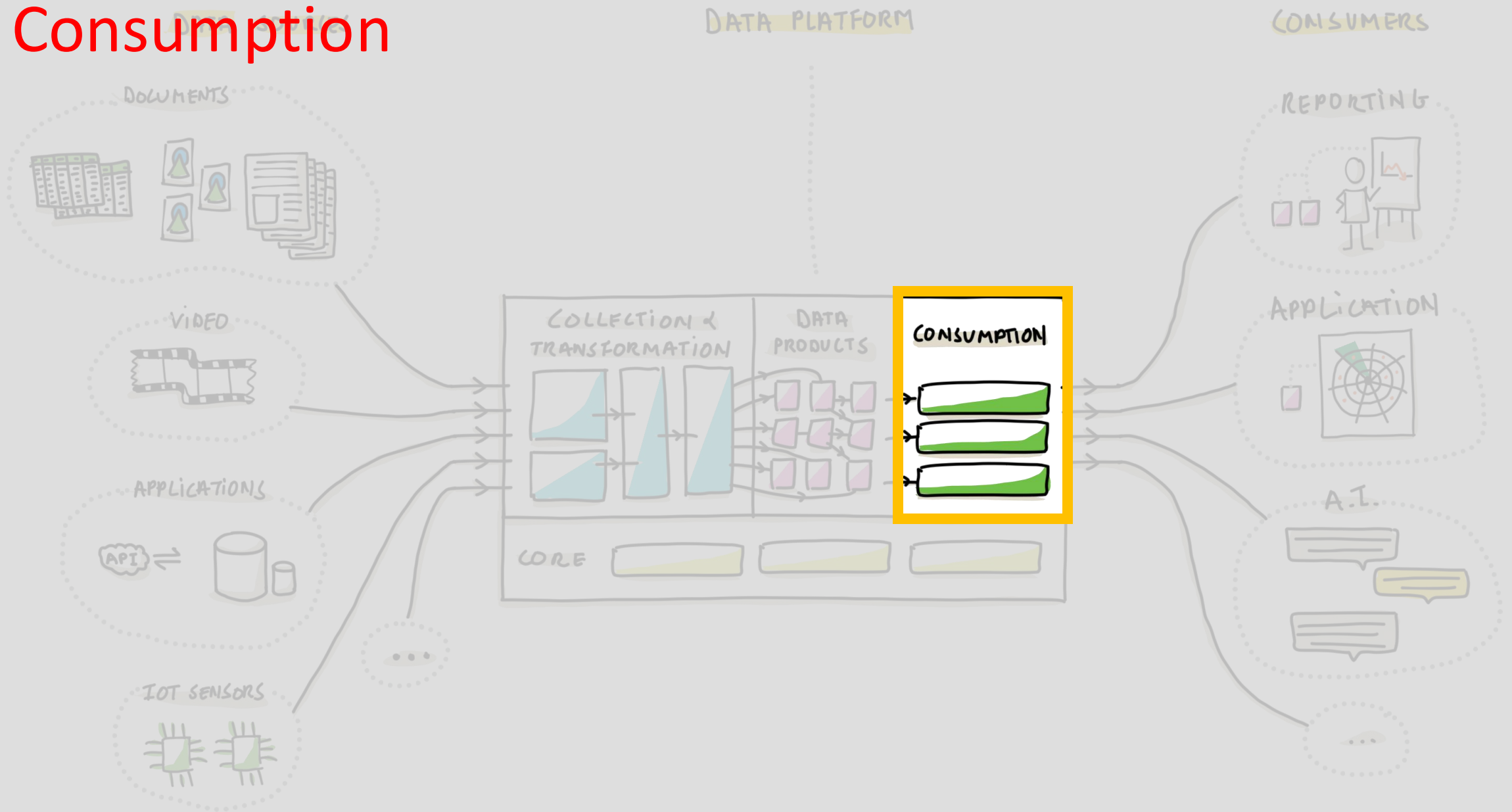
“Operational Plane”

“Analytical Plane”

“Operational / Analytical Plane”



3. Consumption



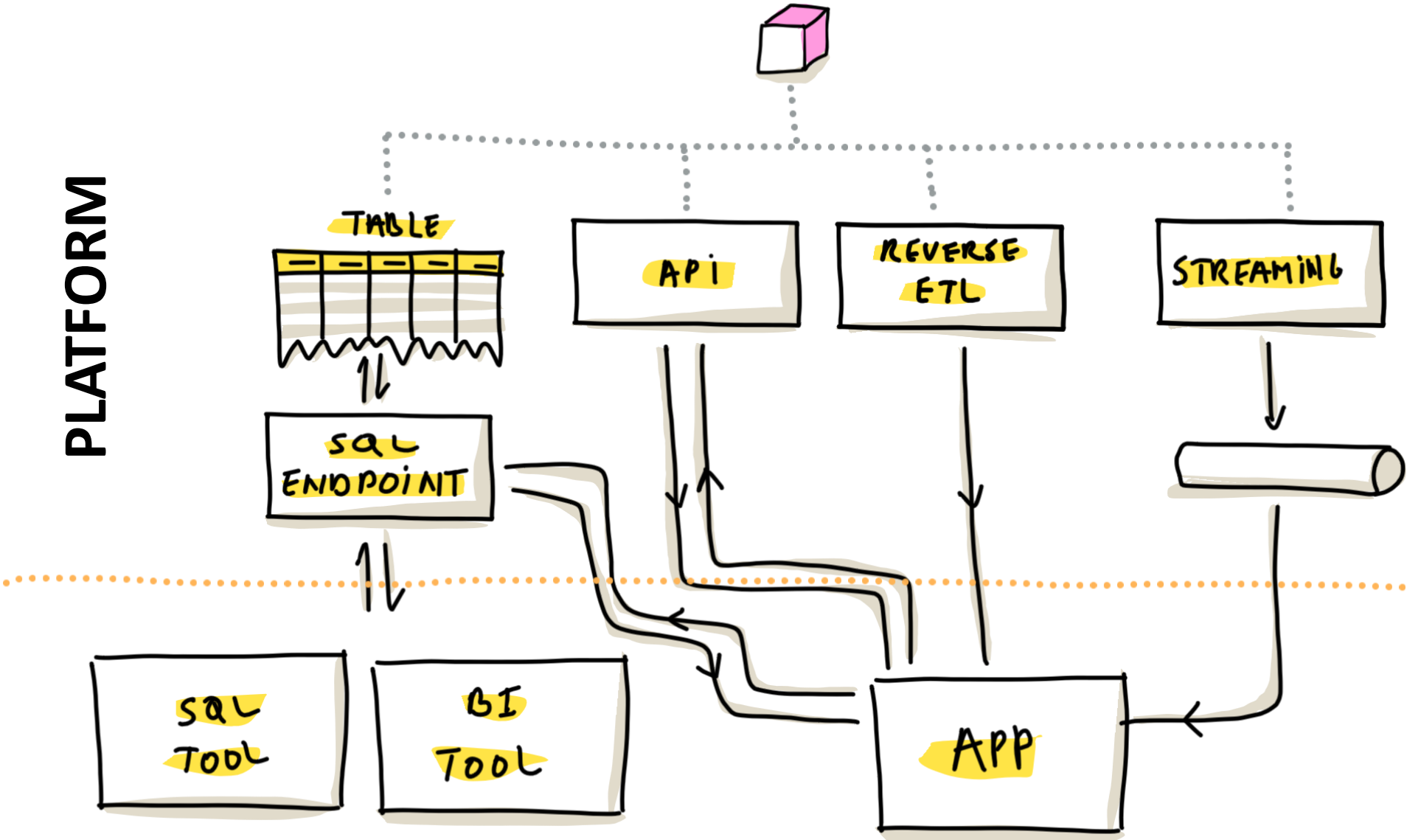
“Operational Plane”

“Analytical Plane”

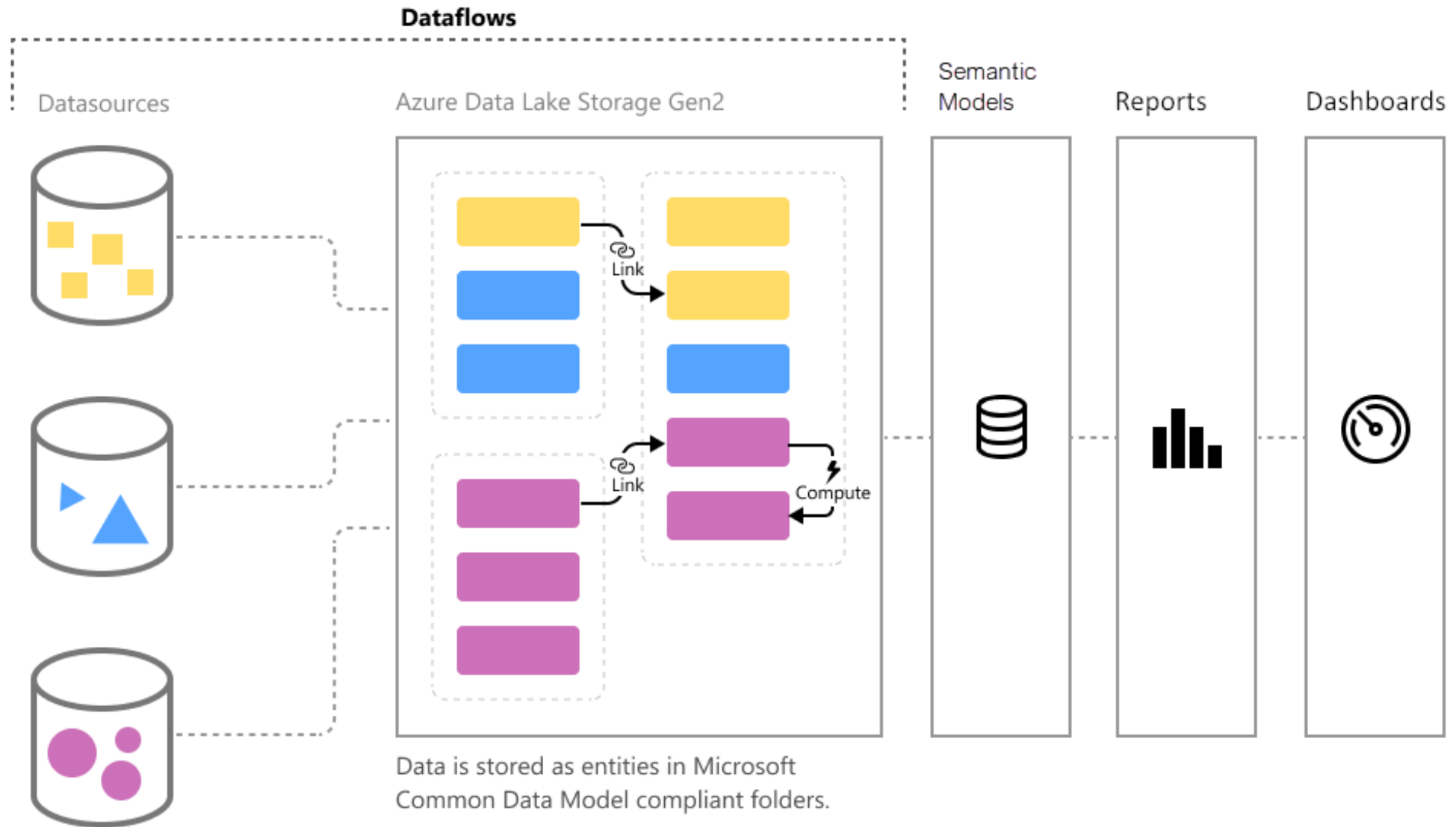
“Operational / Analytical Plane”



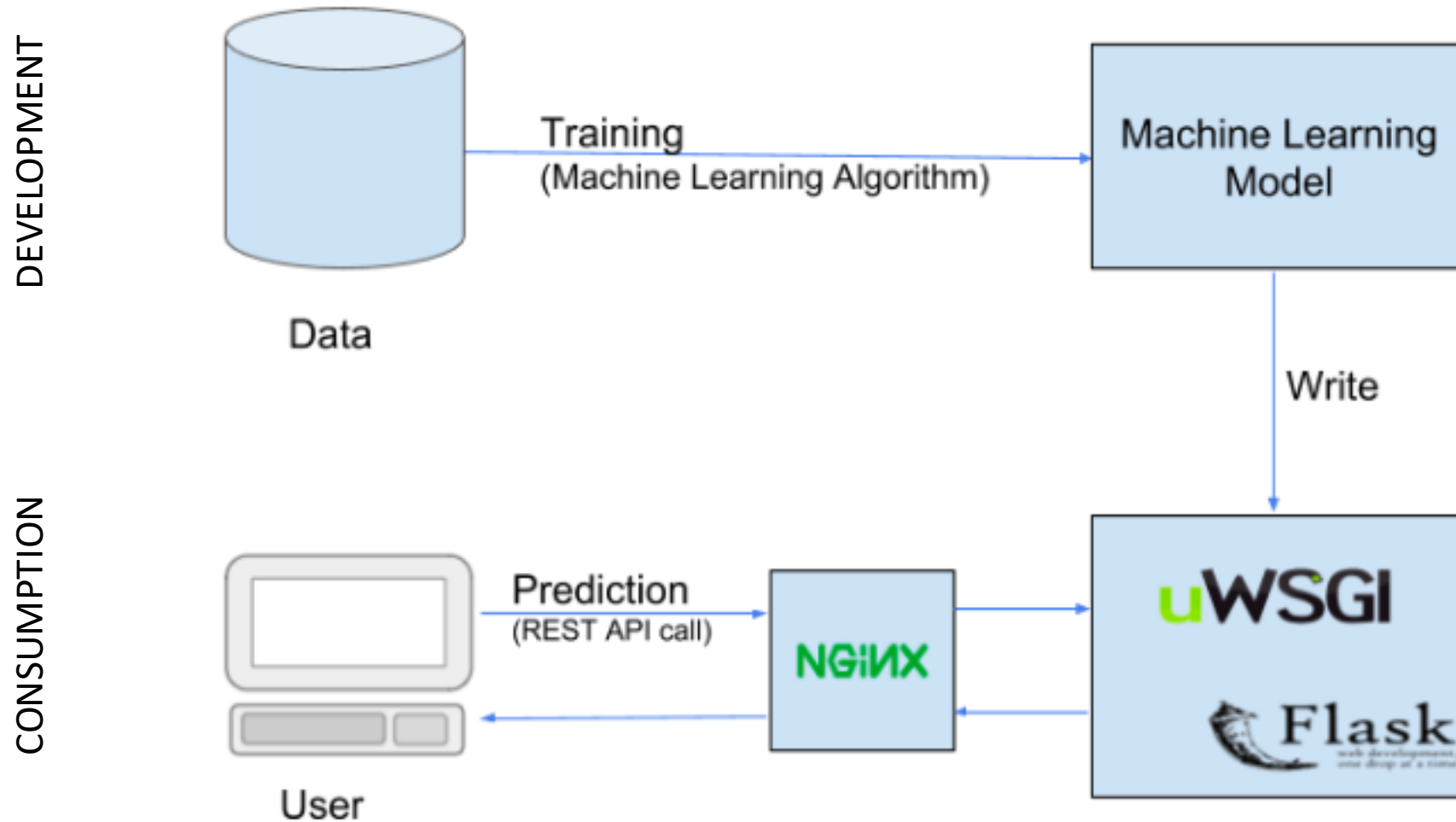
Common Consumption Patterns



Consumption Pattern: Power BI



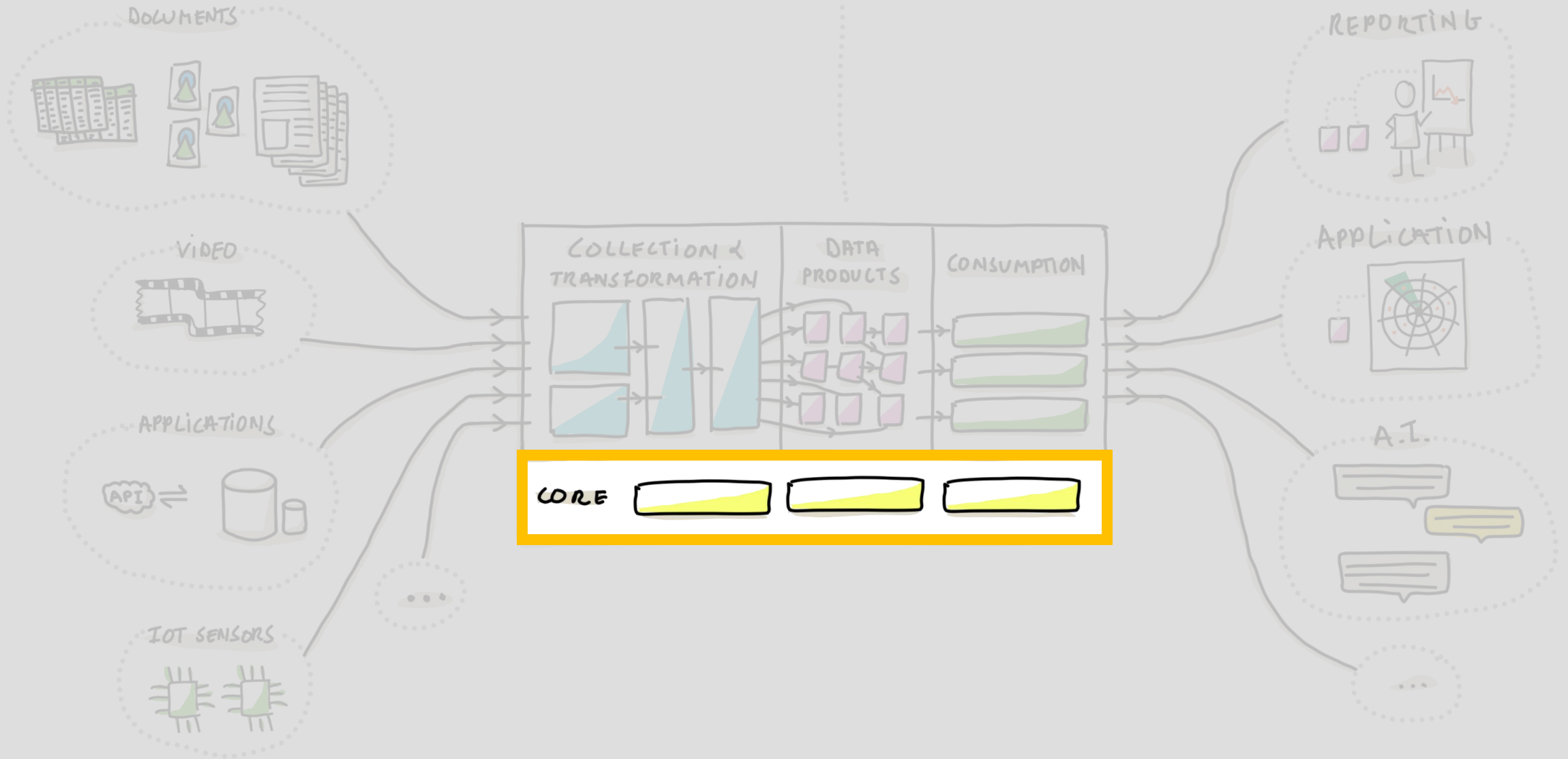
Deploy AI Models as an API



4. Core DATA SOURCES

DATA PLATFORM

CONSUMERS



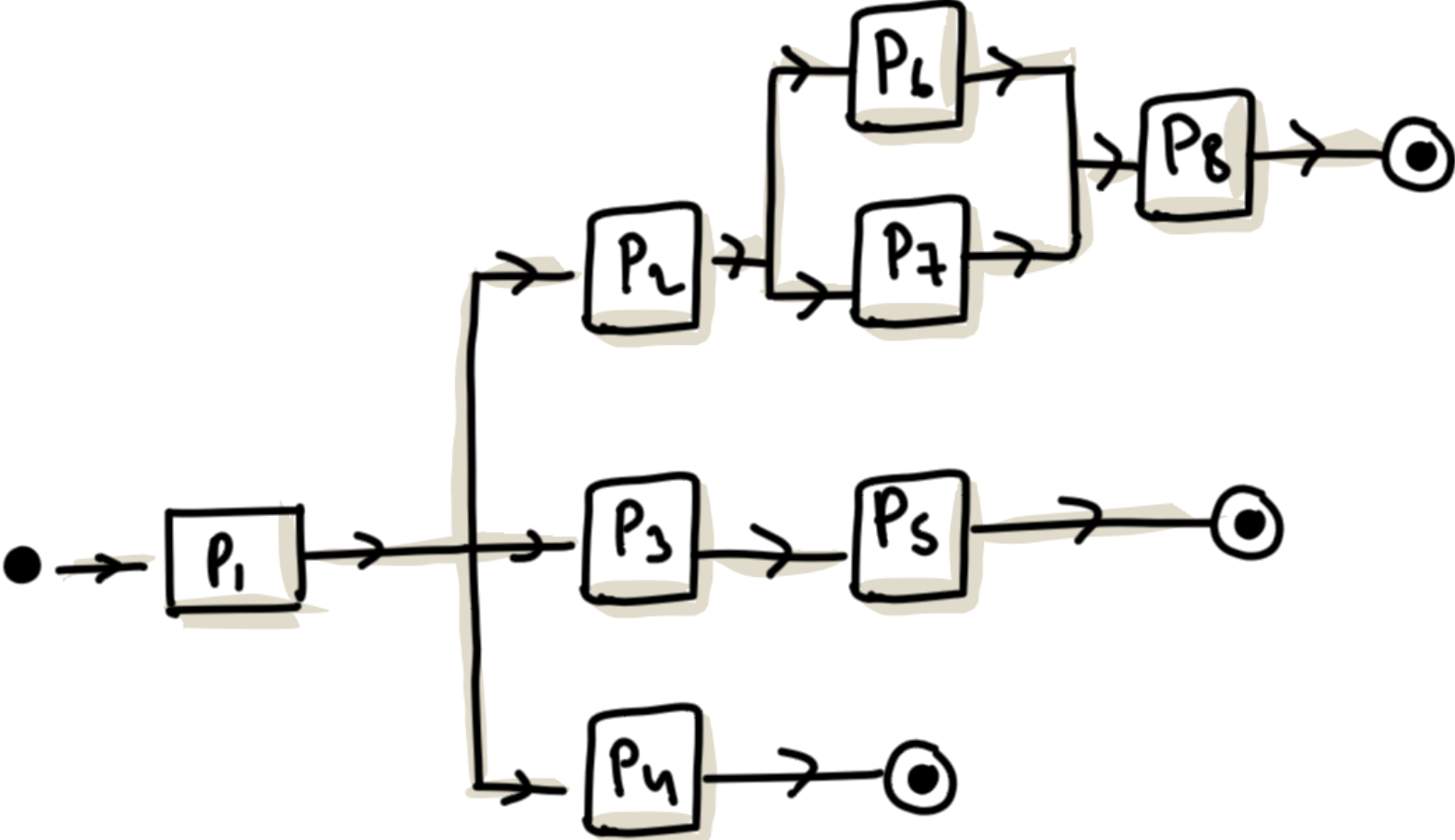
“Operational Plane”

“Analytical Plane”

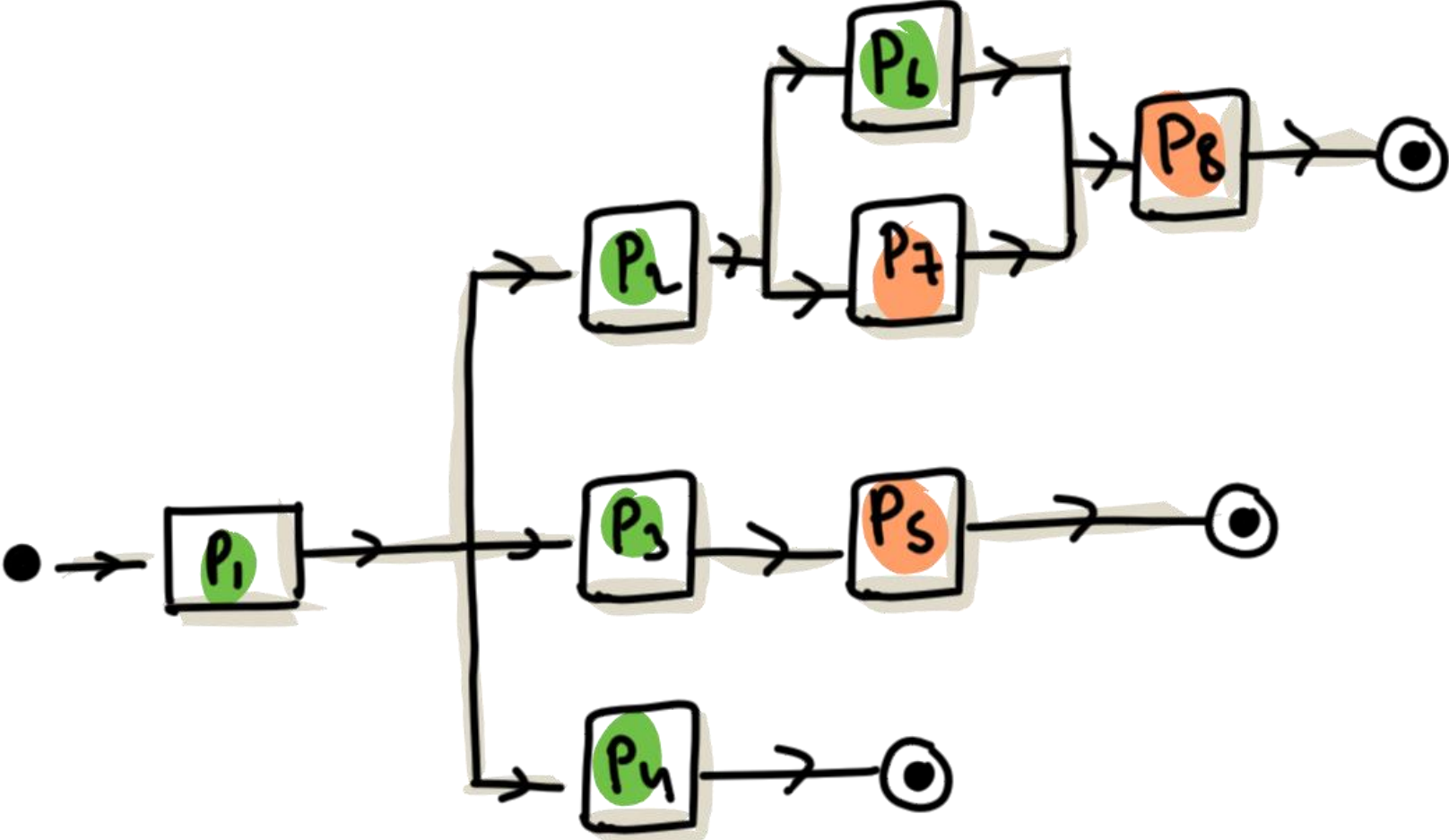
“Operational / Analytical Plane”



Orchestration



Orchestration



Example: Apache Airflow



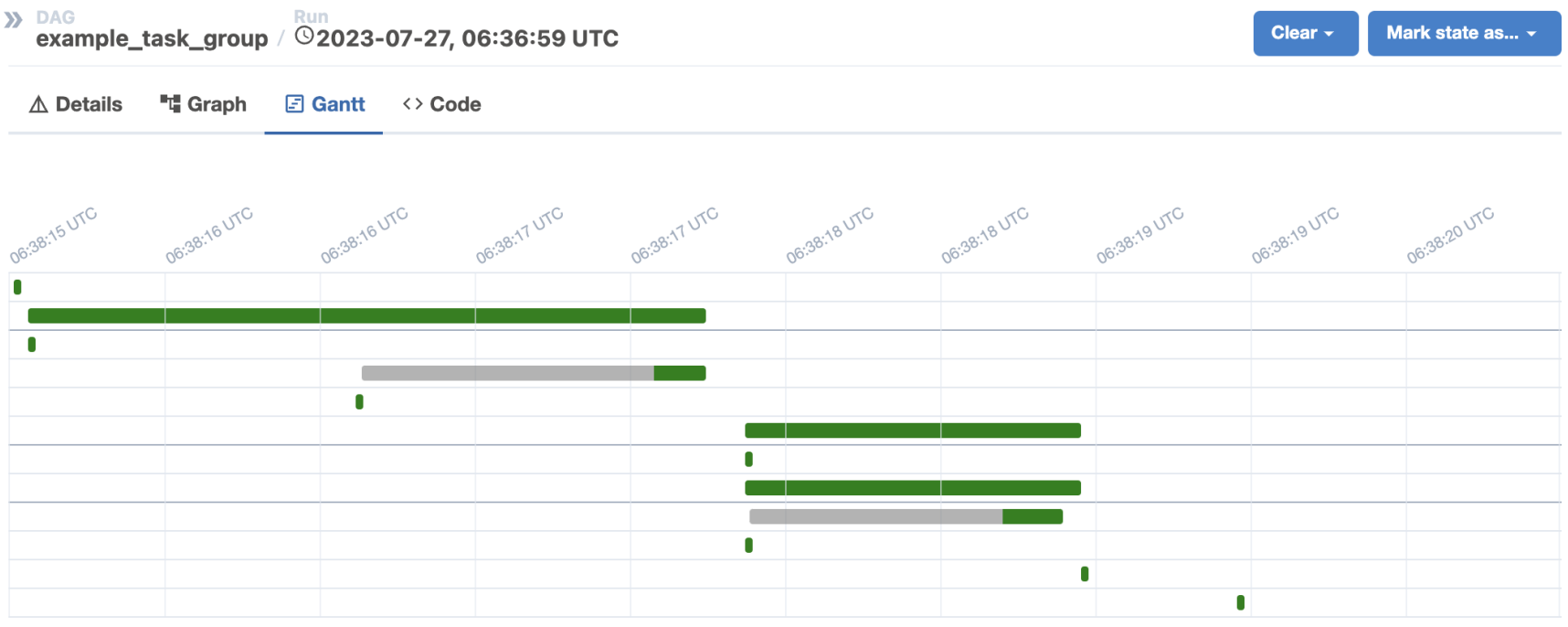
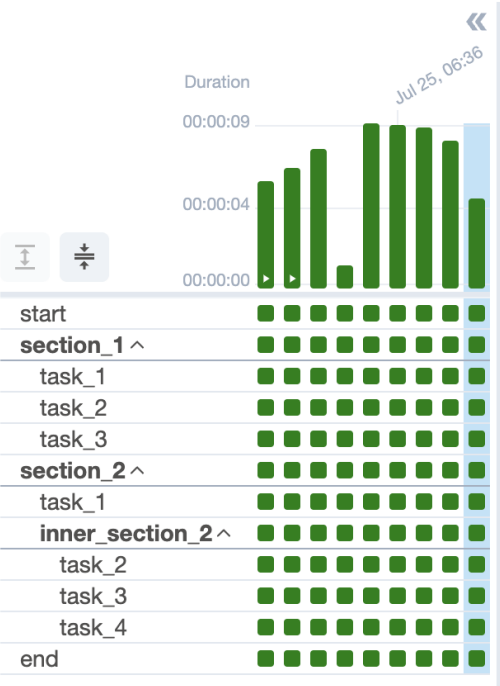
DAG: example_task_group

Schedule: 1 day, 0:00:00 Next Run: 2023-07-27, 06:36:59

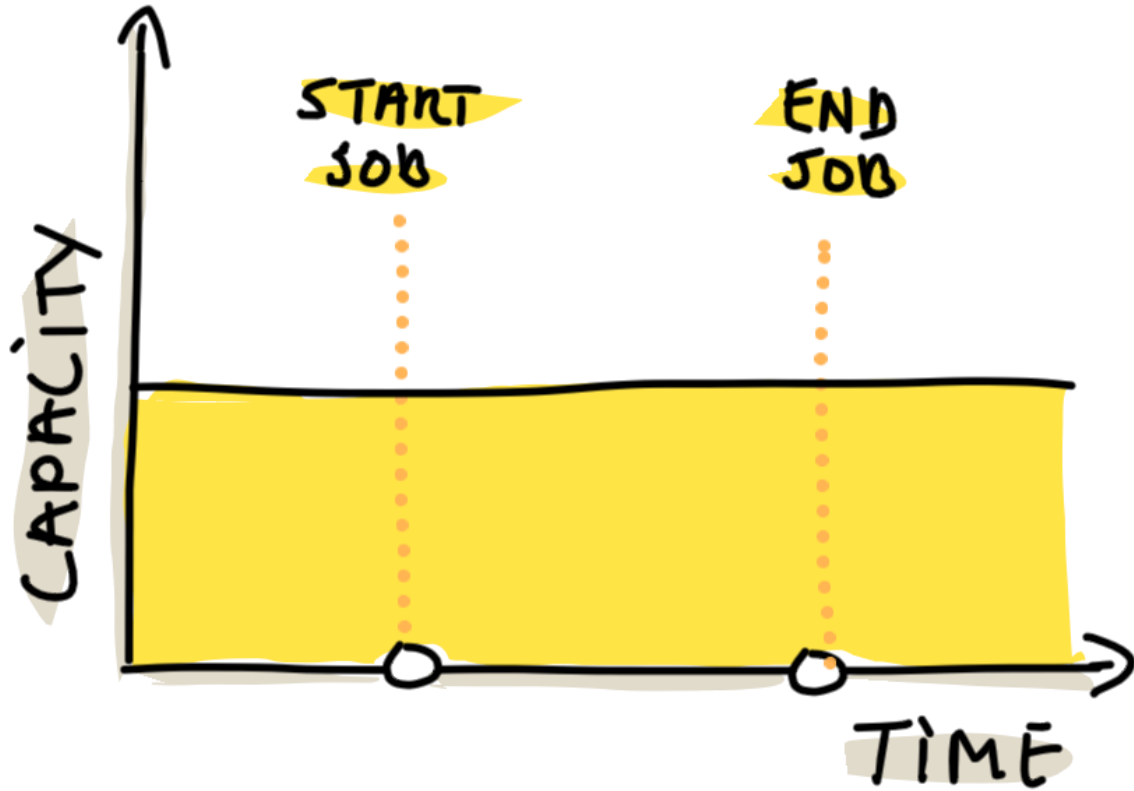
Grid Graph Calendar Task Duration Task Tries Landing Times Gantt Details Code Audit Log

07/28/2023, 06:09:10 AM 25 All Run Types All Run States Clear Filters Auto-refresh

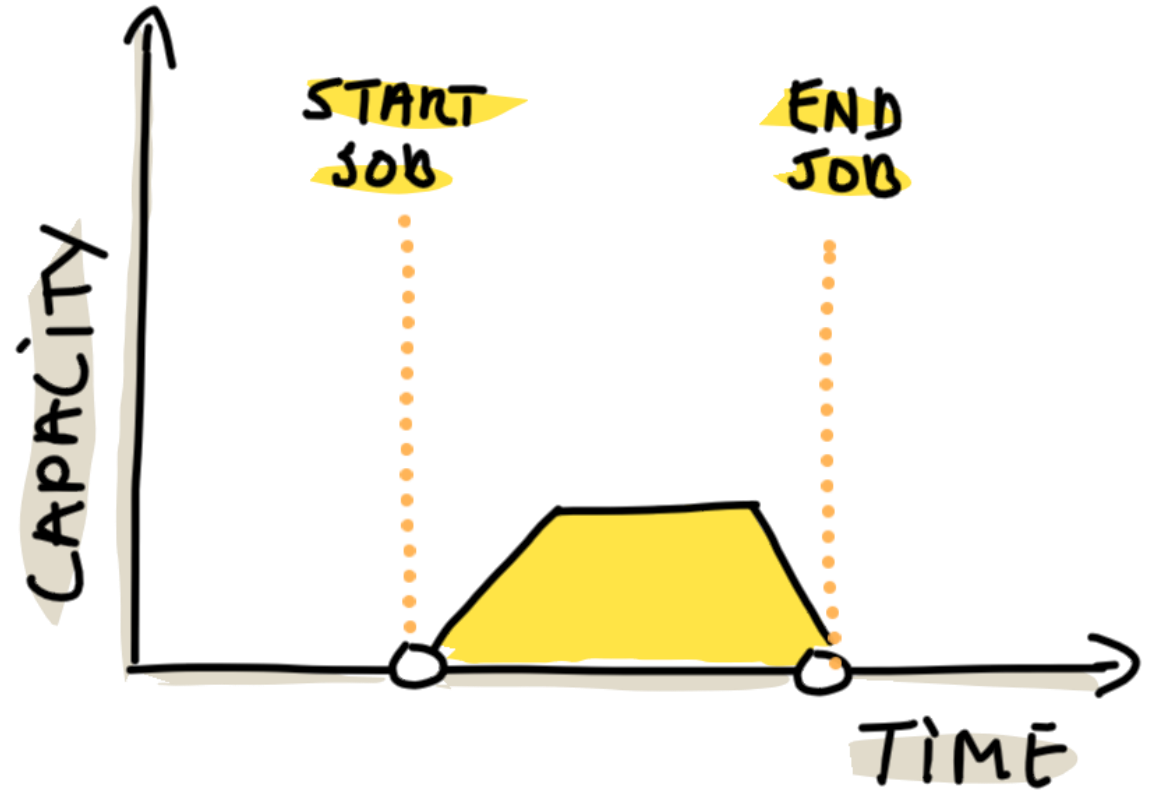
Press shift + / for Shortcuts deferred failed queued removed restarting running scheduled shutdown skipped success up_for_reschedule up_for_retry upstream_failed no_status



Infrastructure Management (Compute)

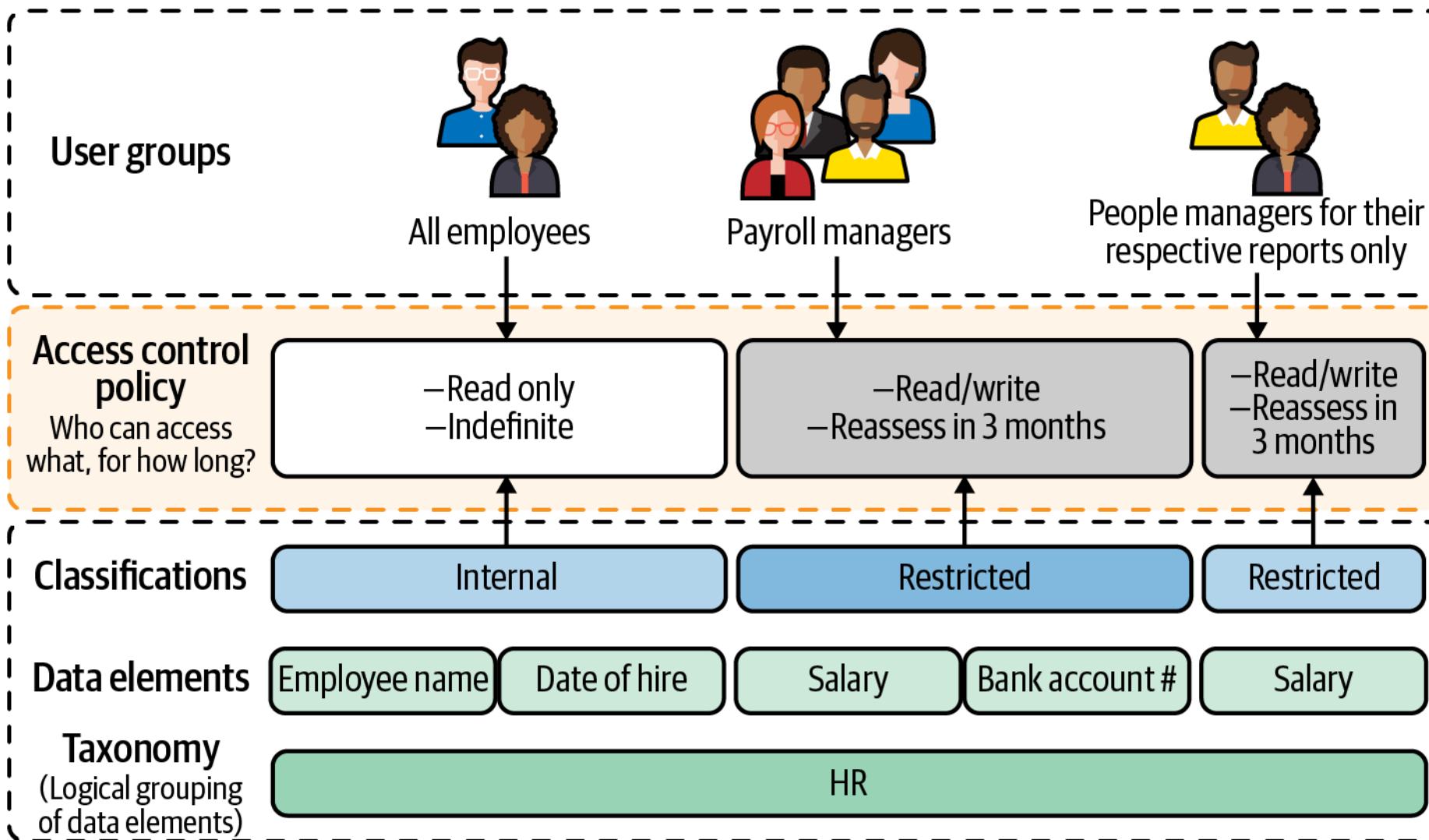


Static



Dynamic

Security & Access Control



Monitoring / Observability

Search anything (⌘J)

TRIAL EXPIRES IN 13 DAYS SCHEDULE CALL

Jerome Williamson will.jerome@green.com

PIPELINES

ACTIVATE

TRANSFORM

DESTINATIONS

Overview

Transformation

Schema Mapper

Load Status

Activity Log

#475 | mysql-source-new MySQL · Ingests every 15 minutes → redshift-destination Amazon Redshift · Loads every 15 minutes

ACTIVE PAUSE

Pipeline Activity

1h 12h 24h ...

Ingestion 08.3K 421.64 epm

Transformations 07.1K 421.64 epm

Schema Mapper 04.2K 421.64 epm

Load 08.3K 421.64 epm

Jobs

Events Ingested

4.8M Events not Loaded

1-10 of 53

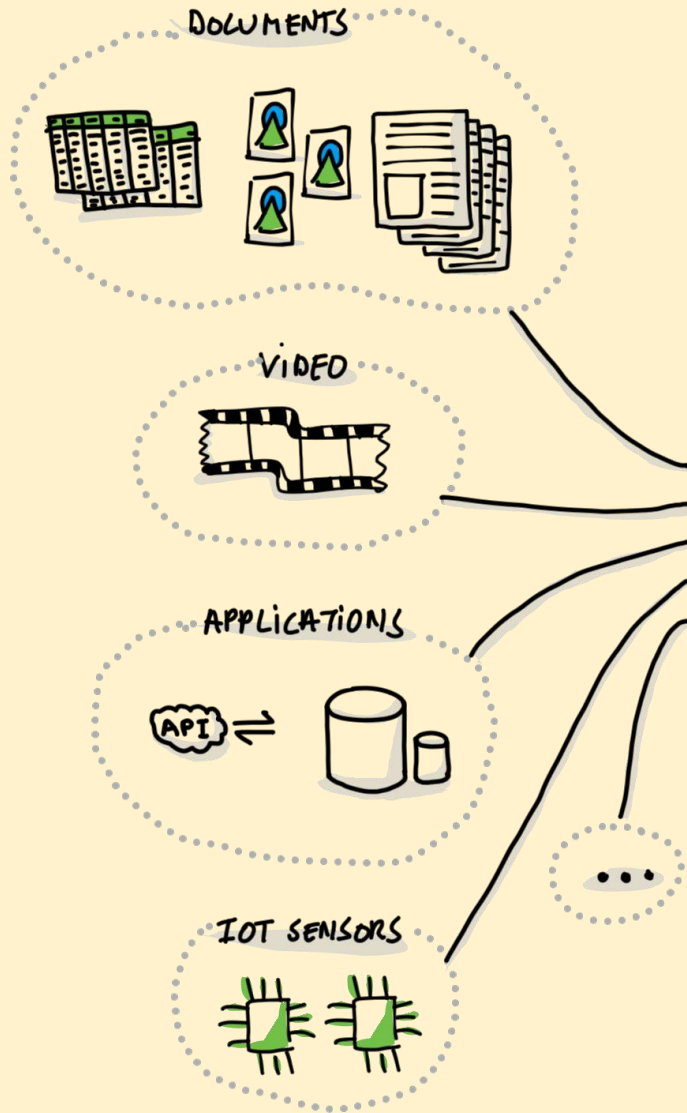
<input type="checkbox"/> object-name_new_updated Historical Load Running	0			QUEUED Not Synced Yet
<input type="checkbox"/> employee_records_updated Position: Sep 26, 2018 3:16:59 PM (UTC)	2.13M		1.28M Events not Loaded	PAUSED Last Synced: 6 Minutes Ago
<input type="checkbox"/> new_customer_data_generated Historical Load Running · Position: Sep 26, 2...	3.43M		3.28M Events not Loaded	ACTIVE Last Synced: 6 Minutes Ago
<input type="checkbox"/> food_categories_new Position: Sep 26, 2018 3:16:59 PM (UTC)	4.21M		18.29K Events not Loaded	FAILED Last Synced: 6 Minutes Ago
<input type="checkbox"/> register_file_loads Position: Sep 26, 2018 3:16:59 PM (UTC)	2.43M		1.38M Events not Loaded	ACTIVE Last Synced: 6 Minutes Ago
<input type="checkbox"/> new_customer_data_generated				ACTIVE

DOCS

LIVE CHAT

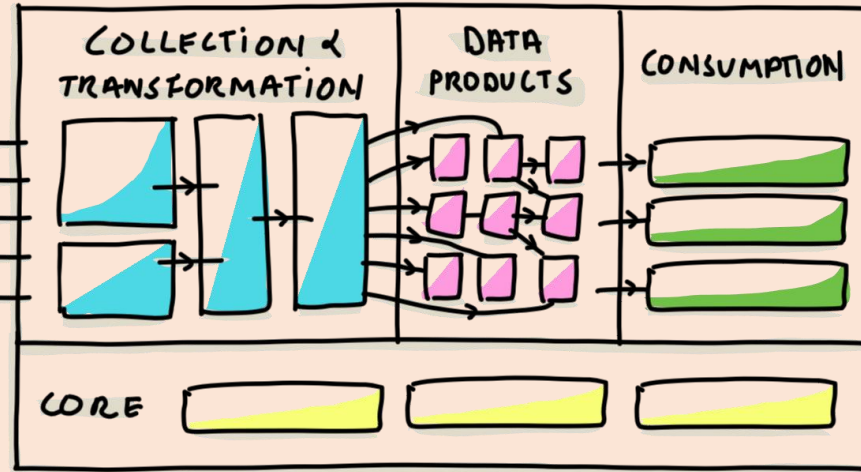
v1.38

DATA SOURCES



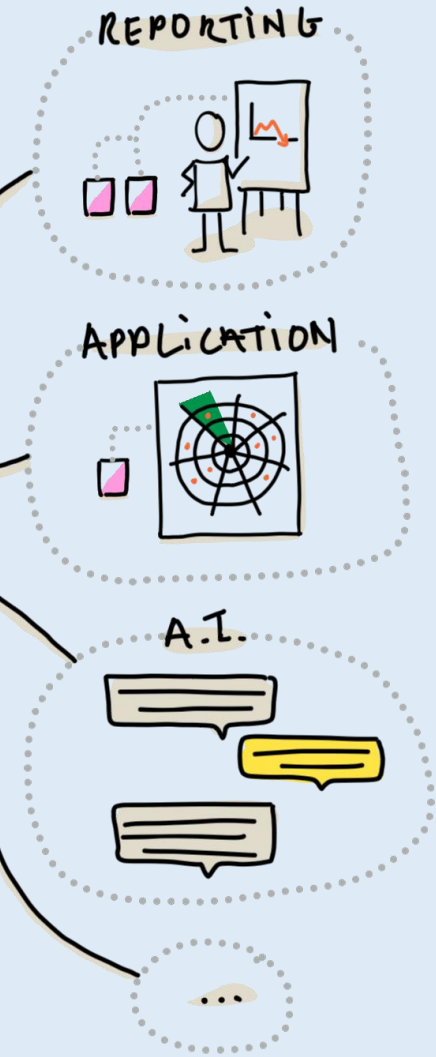
“Operational Plane”

DATA PLATFORM



“Analytical Plane”

CONSUMERS



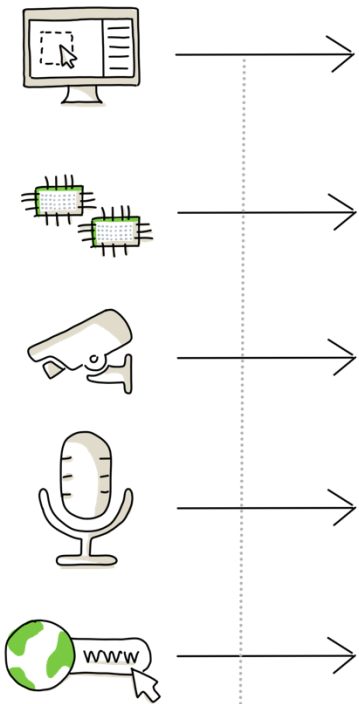
“Operational / Analytical Plane”

Table of Contents

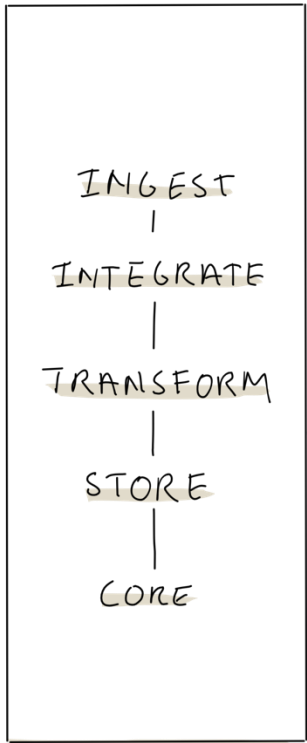
- Dead Horse Theory
- Data Platform
 - Introduction
 - Core Layers
 - **Additional Layers**
- Technology Selection



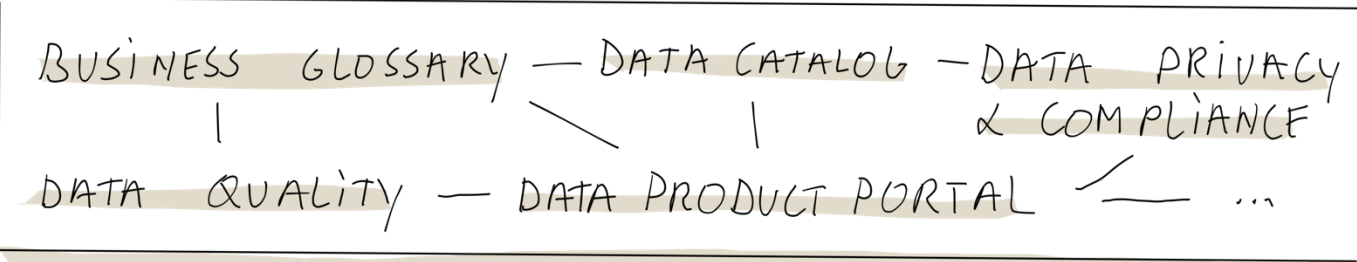
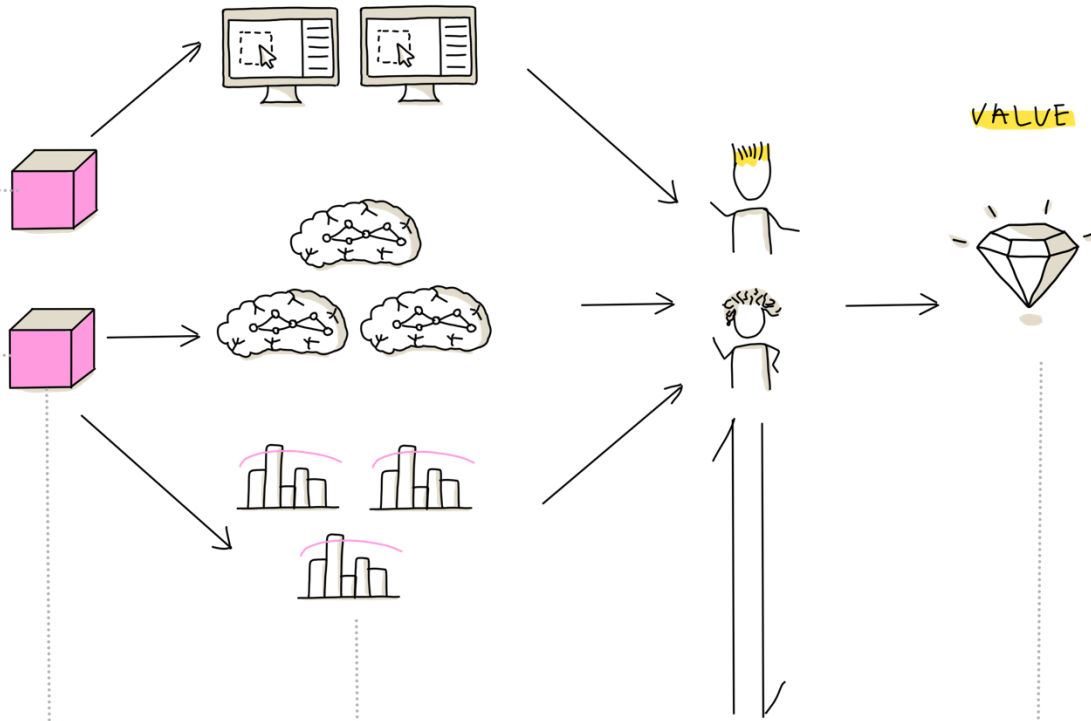
DATA PRODUCERS (SOURCES)



DATA PLATFORM



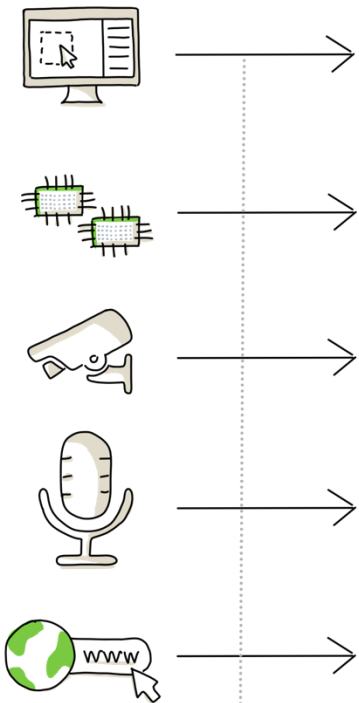
CONSUMPTION



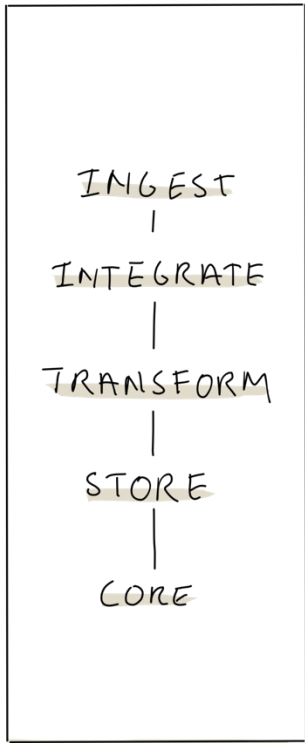
META-DATA MANAGEMENT



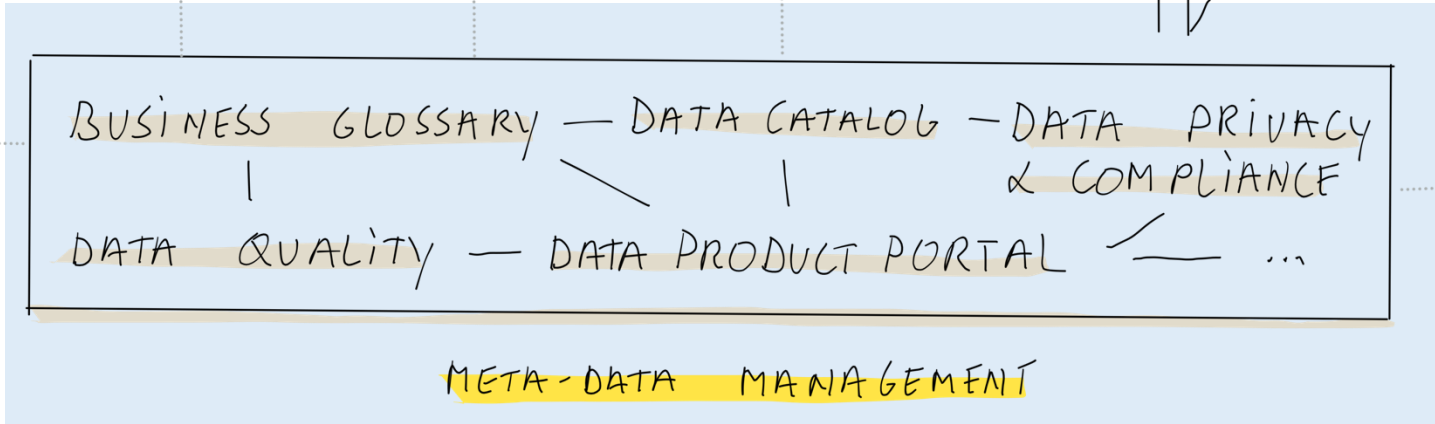
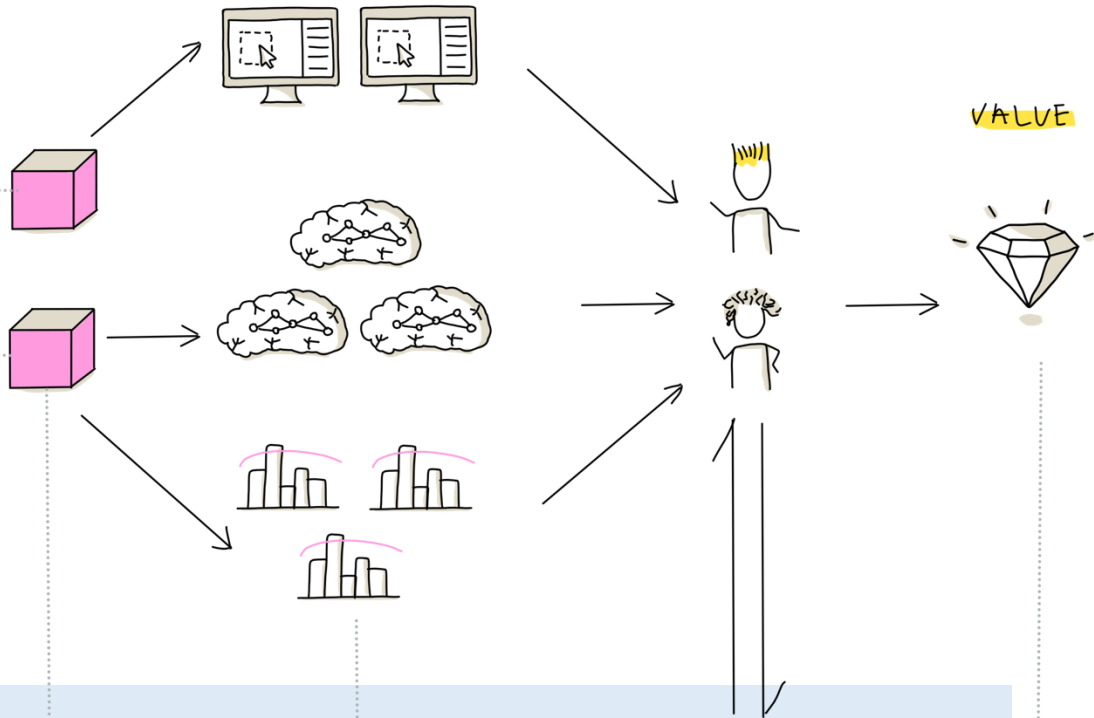
DATA PRODUCERS (SOURCES)



DATA PLATFORM



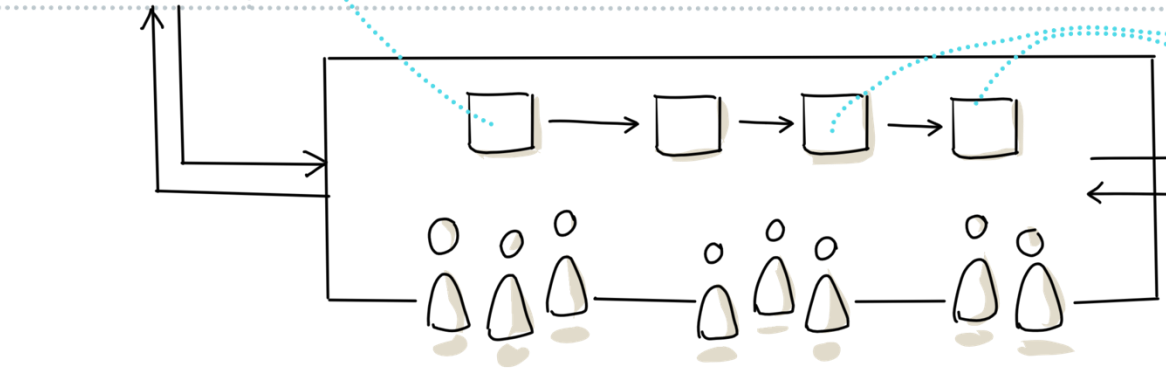
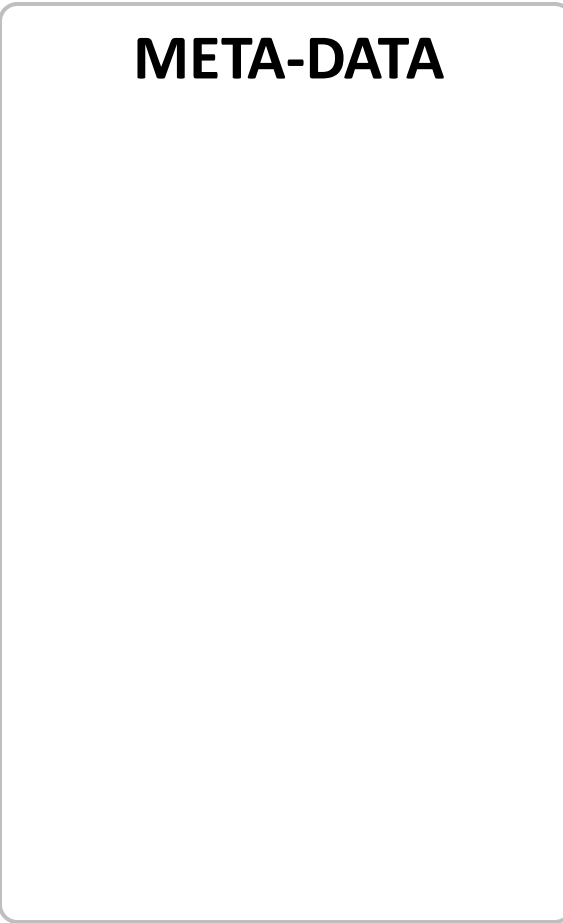
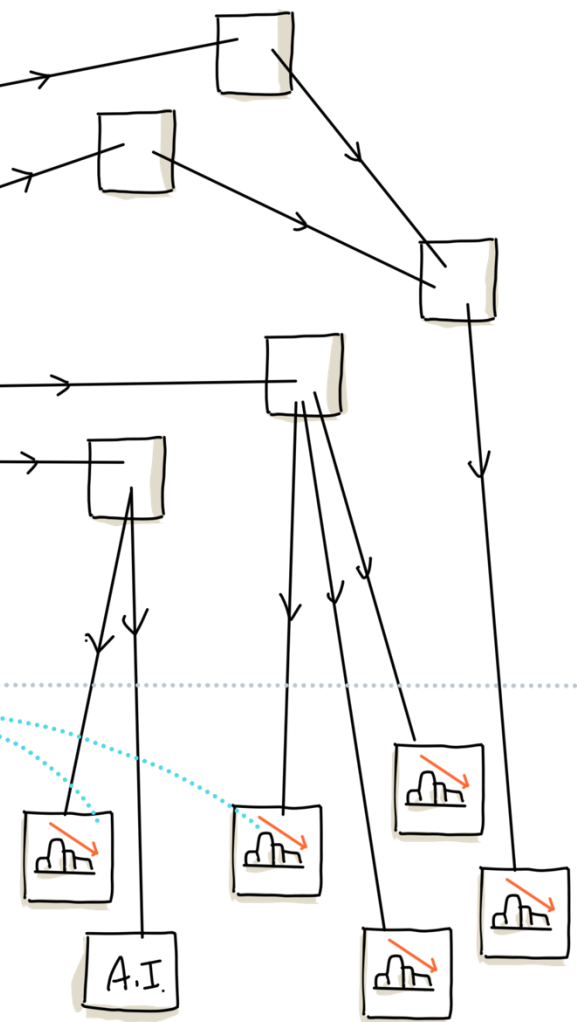
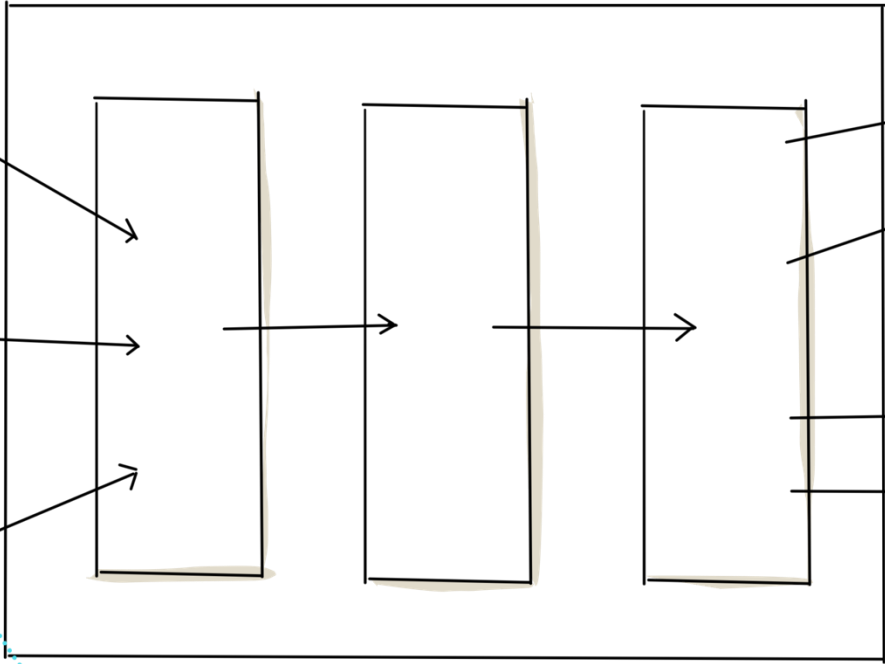
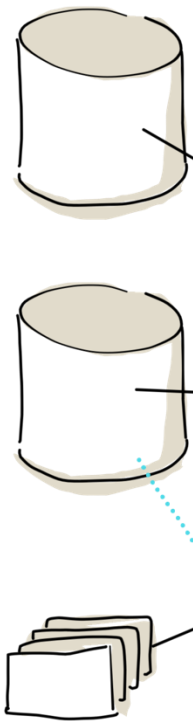
CONSUMPTION



DATA SOURCES

DATA TRANSFORMATIONS

DATA PRODUCTS/
SETS



USERS

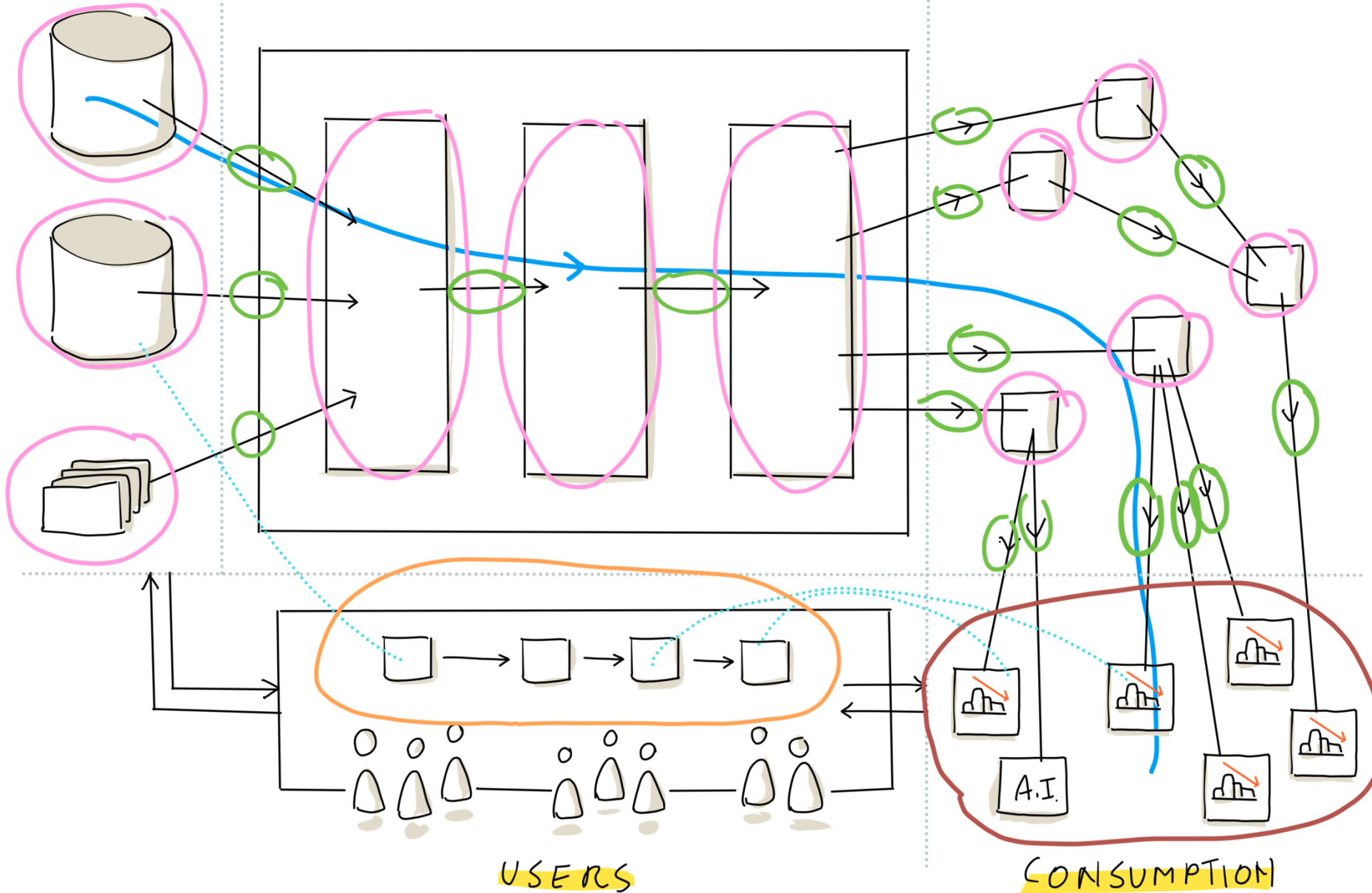
CONSUMPTION



DATA SOURCES

DATA TRANSFORMATIONS

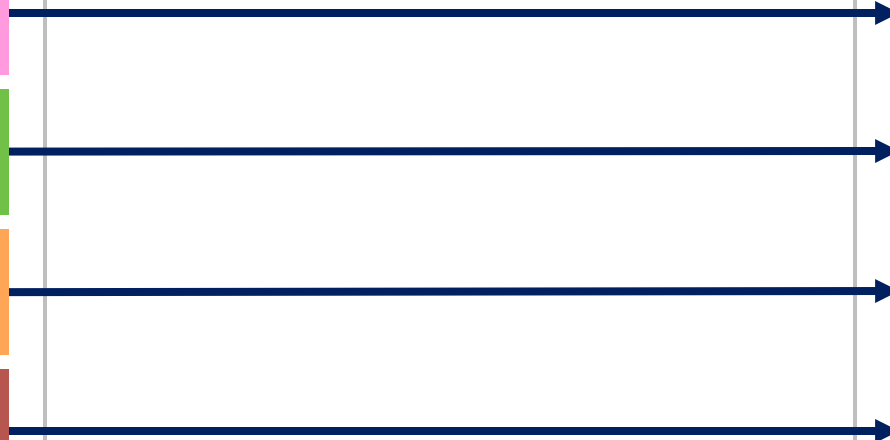
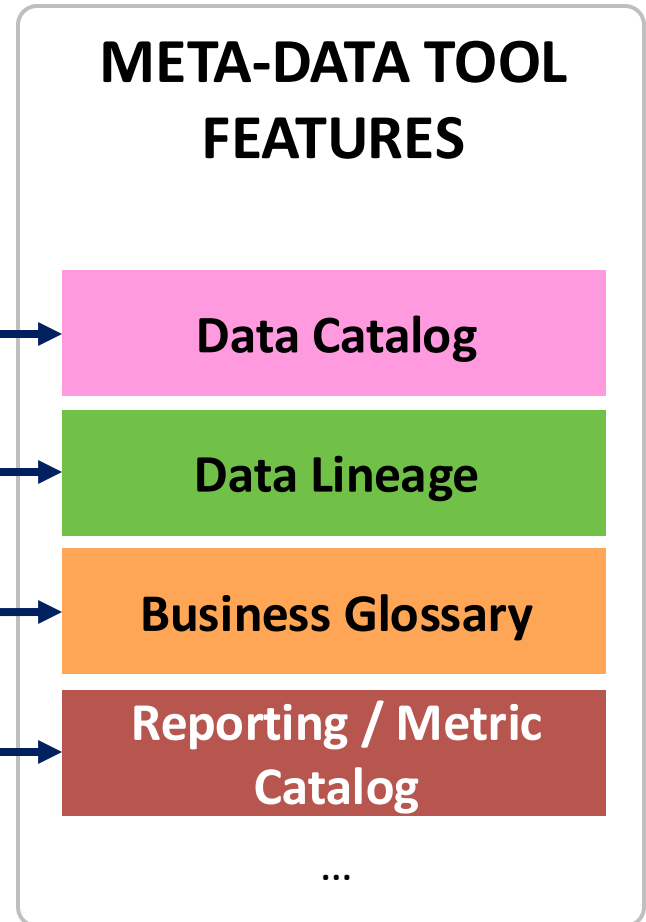
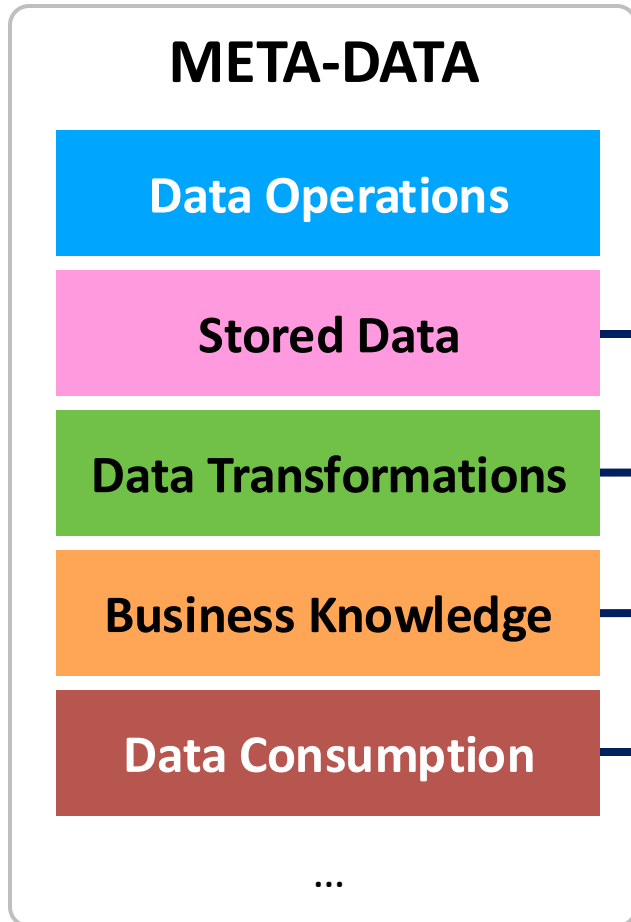
DATA PRODUCTS/
SETS



META-DATA

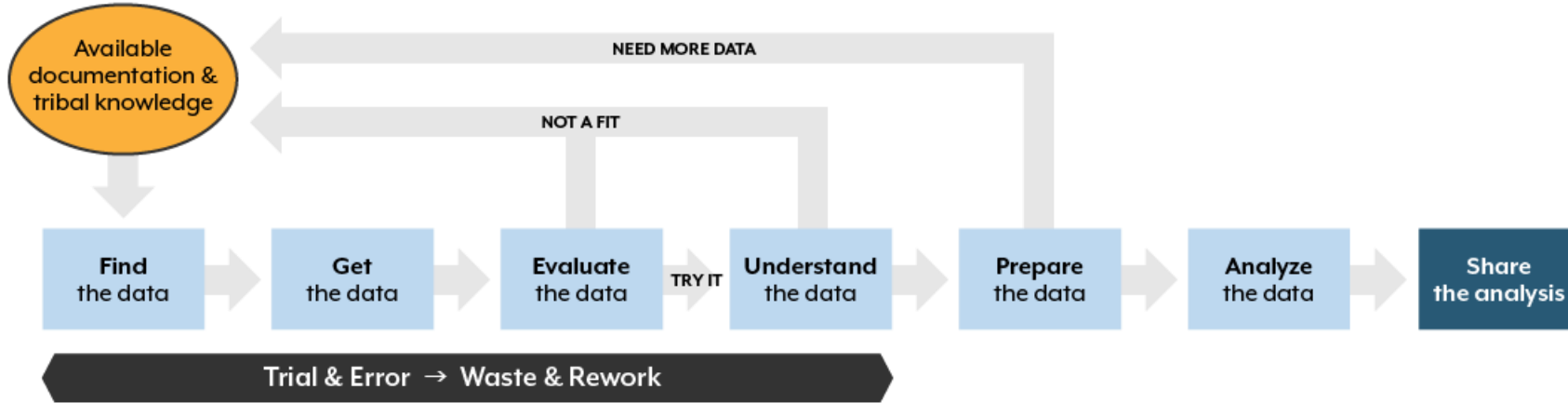
- Data Operations
- Stored Data
- Data Transformations
- Business Knowledge
- Data Consumption
- ...



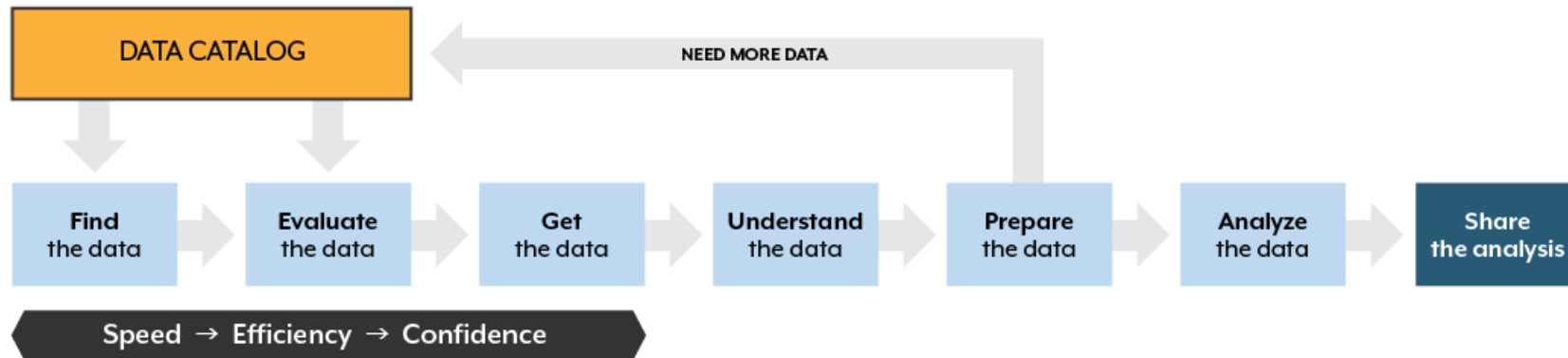


DATA CATALOG

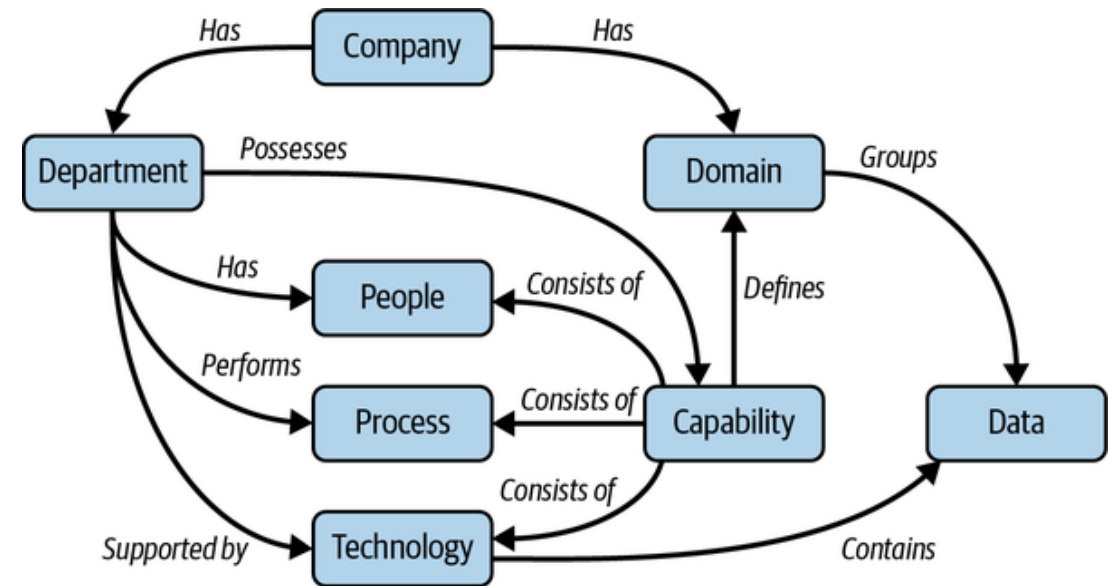
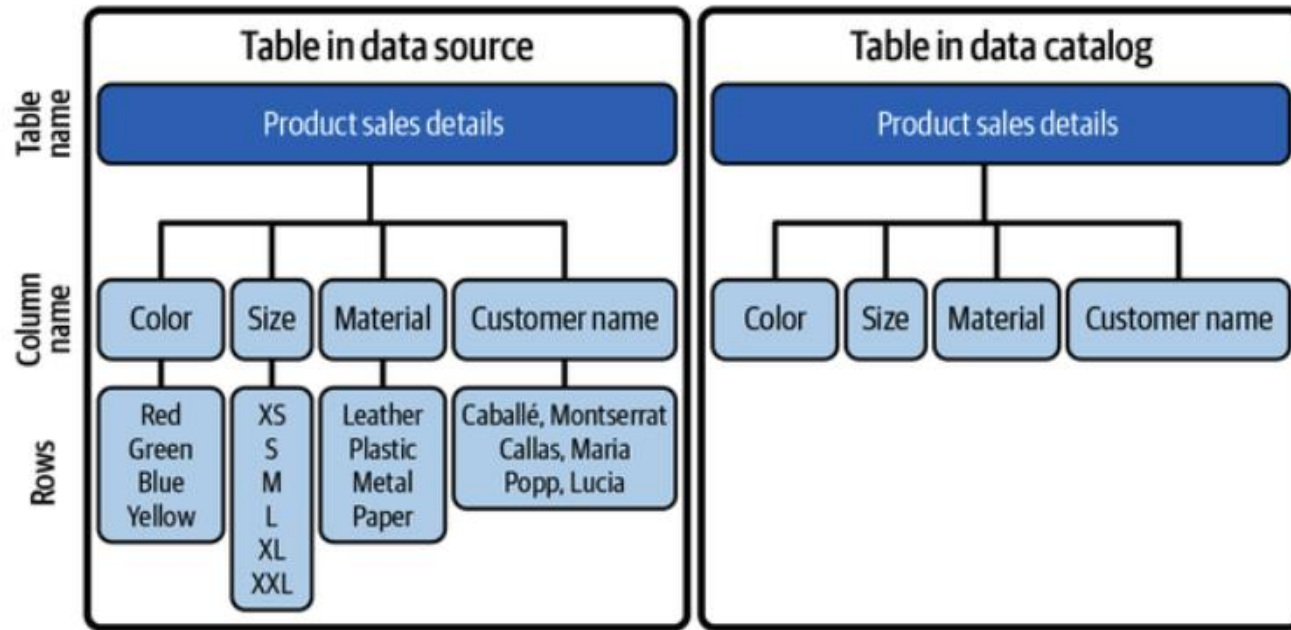
Without Data Catalog



With Data Catalog



DATA CATALOG



DATA CATALOG



View all

36 K

Analytics

Govern



Schedule a Demo

- Platform
- Domain
- Type
- Glossary Term
- Owned By
- + More Filters

Advanced Filters

Navigate

Showing 1 - 10 of 44 results

- Charts 5
- Dashboards 2
- Datasets 36
- Pipelines 0

Dataset PostgreSQL

customers

This table has basic information about a customer, as well as some derived facts based on a customer's orders

Matches column customer_id Matches column description Customer's first n...

Table dbt & BigQuery calm-pagoda-323403 jaffle_shop

customers

This table has basic information about a customer, as well as some derived facts based on a customer's orders

Bronze dbt:contains_pii

Matches column customer_id Matches column description Customer's first n...

customers customers

Owners

bi-engineering influencer

Table BigQuery calm-pagoda-323403 jaffle_shop

customers_source

Dataset PostgreSQL

customers_source

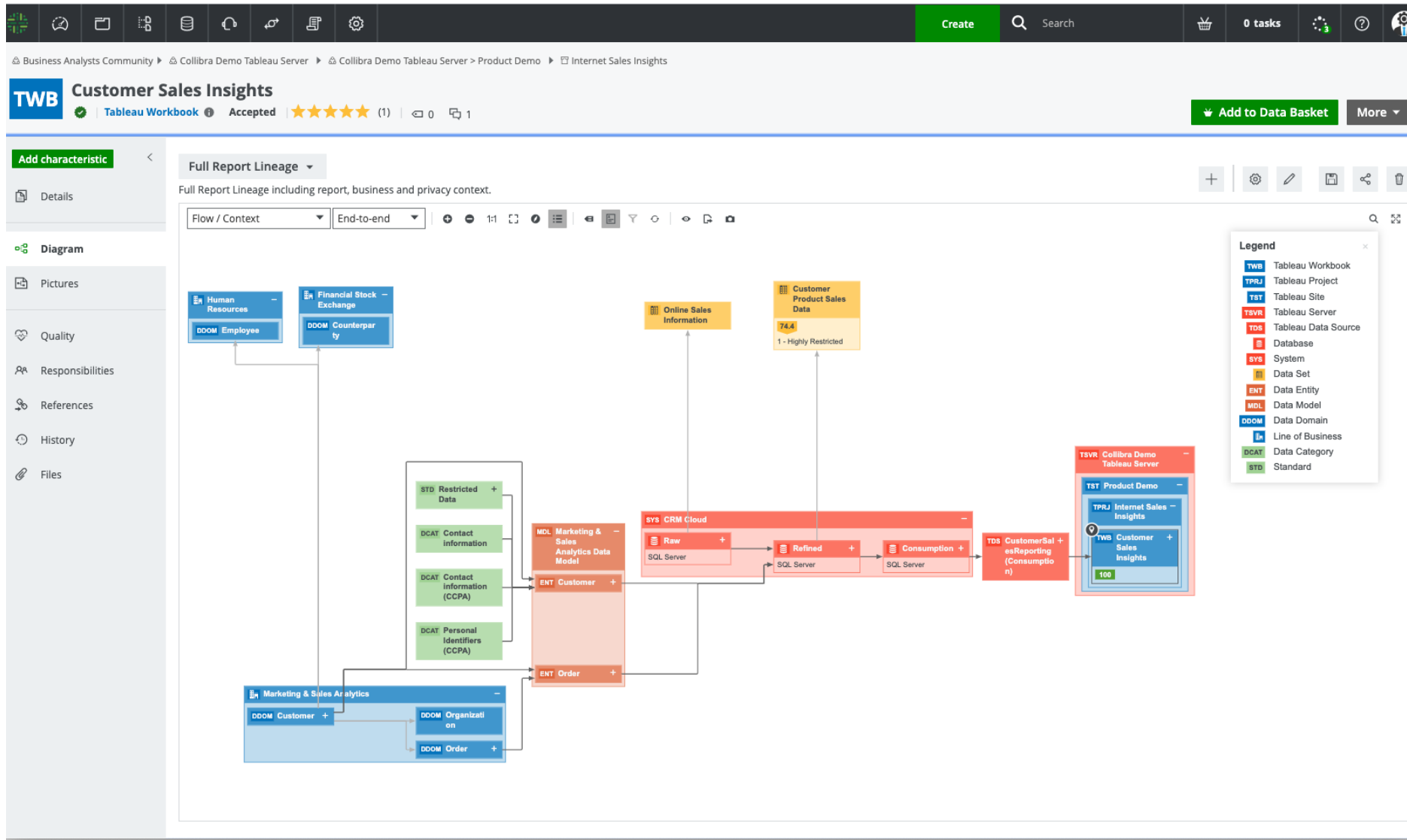
Table dbt & Snowflake long_tail_companions analytics

customer_last_purchase_date

Owners



DATA LINEAGE



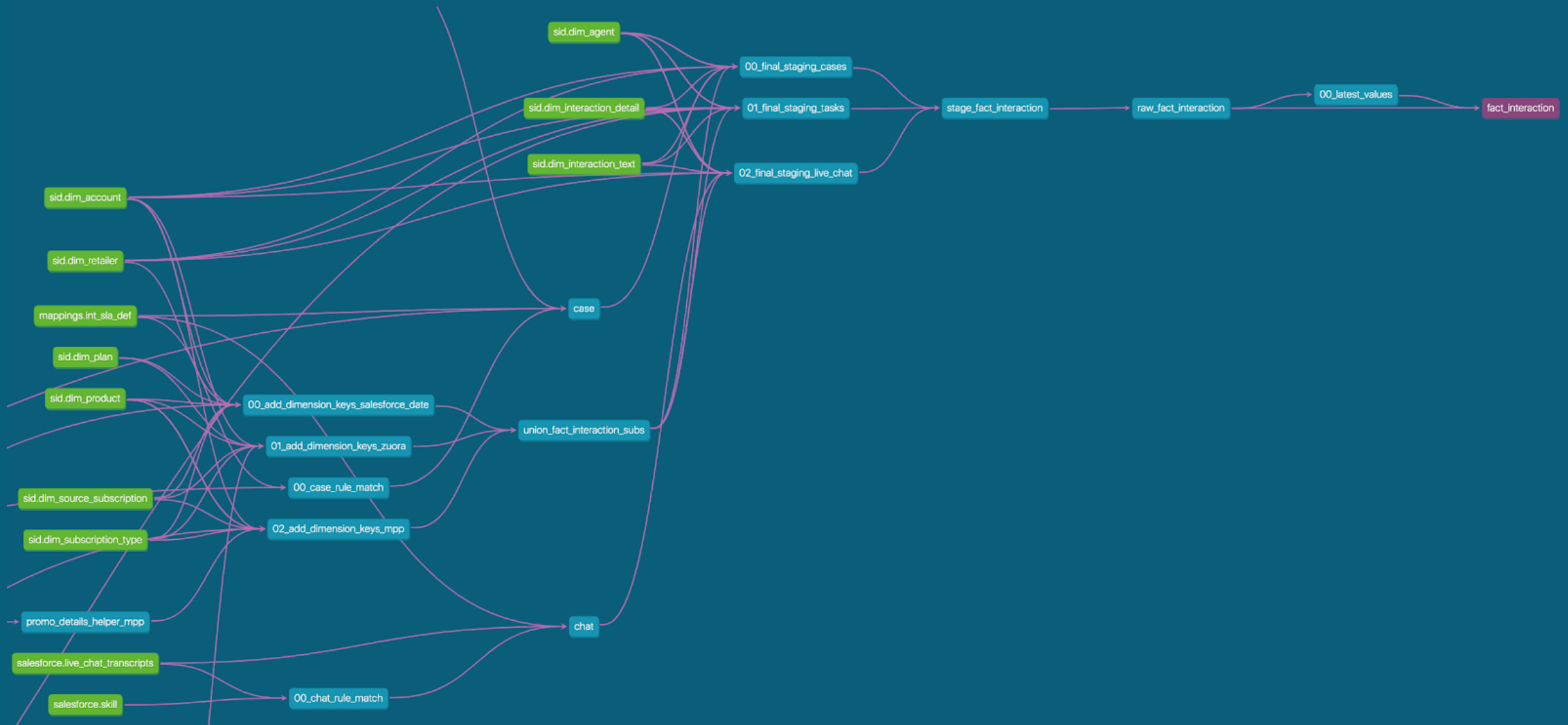
[Data Lineage in Colibra]



DATA LINEAGE



DATA LINEAGE



BUSINESS GLOSSARY

The screenshot displays the Atlan Business Glossary interface. On the left, a sidebar contains navigation options: 'Assets', 'Glossary', and 'Insights'. The 'Glossary' section is expanded, showing a search bar and a list of categories including 'Aisle', 'Concepts', 'Consumer Product Goods', 'COVID-19', 'Example Glossary', 'Instacart', 'KPIs', and 'Metrics'. Under 'Metrics', 'Customer Acquisition Cost' is selected and highlighted.

The main content area shows the 'Customer Acquisition Cost' term page. At the top, there are navigation icons, a title 'Customer Acquisition Cost', and a 'TERM' icon. Below the title, there are tabs for 'Overview' and 'Linked Assets'. The 'Overview' tab is active, displaying a 'Readme' section with an 'Edit' button.

The 'Simple method' section explains that the simple method divides the total marketing costs to acquire new customers by the total number of customers acquired in a defined period. It includes the formula:

$$CAC = \frac{MCC}{CA}$$

where:

- CAC = Customer Acquisition Cost
- MCC = total marketing cost for acquiring customers (not regular customers)
- CA = total customers acquired

The 'Complex method' section explains that in addition to the costs incurred in marketing, the complex method includes sales and marketing wages, software costs for sales and marketing, all additional professional services such as designers, consultants, etc., as well as other overhead costs. It includes the formula:

$$CAC = \frac{MCC + W + S + PS + O}{CA}$$

where:

- CAC = Customer Acquisition Cost
- MCC = total marketing cost for acquiring customers (not regular customers)
- W = wages connected with sales and marketing
- S = all the marketing and sales associated software cost (inc. E-Commerce-Platform,

On the right side of the interface, there is an 'Overview' panel with various metadata fields: 'Owners' (chris), 'Classification' (Confidential), 'Certificate' (Verified, chris 3 months ago), 'Categories' (This term does not belong to any category), 'Related Terms' (Average Selling Price, Churn Rate, Customer Lifetime Value), and 'Custom Metadata' (Great Expectations, PO number, Airflow ETL Details). A vertical sidebar on the far right contains additional navigation options like 'Activity', 'Resources', 'Request', 'Property', 'Great...', 'Data F...', 'PO number', 'Data Q...', 'Data C...', 'Airflo...', and 'Priority'.



REPORTING / METRIC CATALOG

The screenshot displays the Tableau Reports interface. At the top, there is a navigation bar with icons for home, reports, data sets, data sources, data dictionary, technology assets, metrics, access requests, and advanced data types. A 'Create' button and a search bar are also present. Below the navigation bar, the main content area is titled 'Tableau Reports' and includes a 'Revert to original' button with a timestamp 'last changes 2 minutes ago'. The interface is organized into a grid of report cards, each featuring a title, a 'Candidate' status, and a small preview of the report content. The reports include:

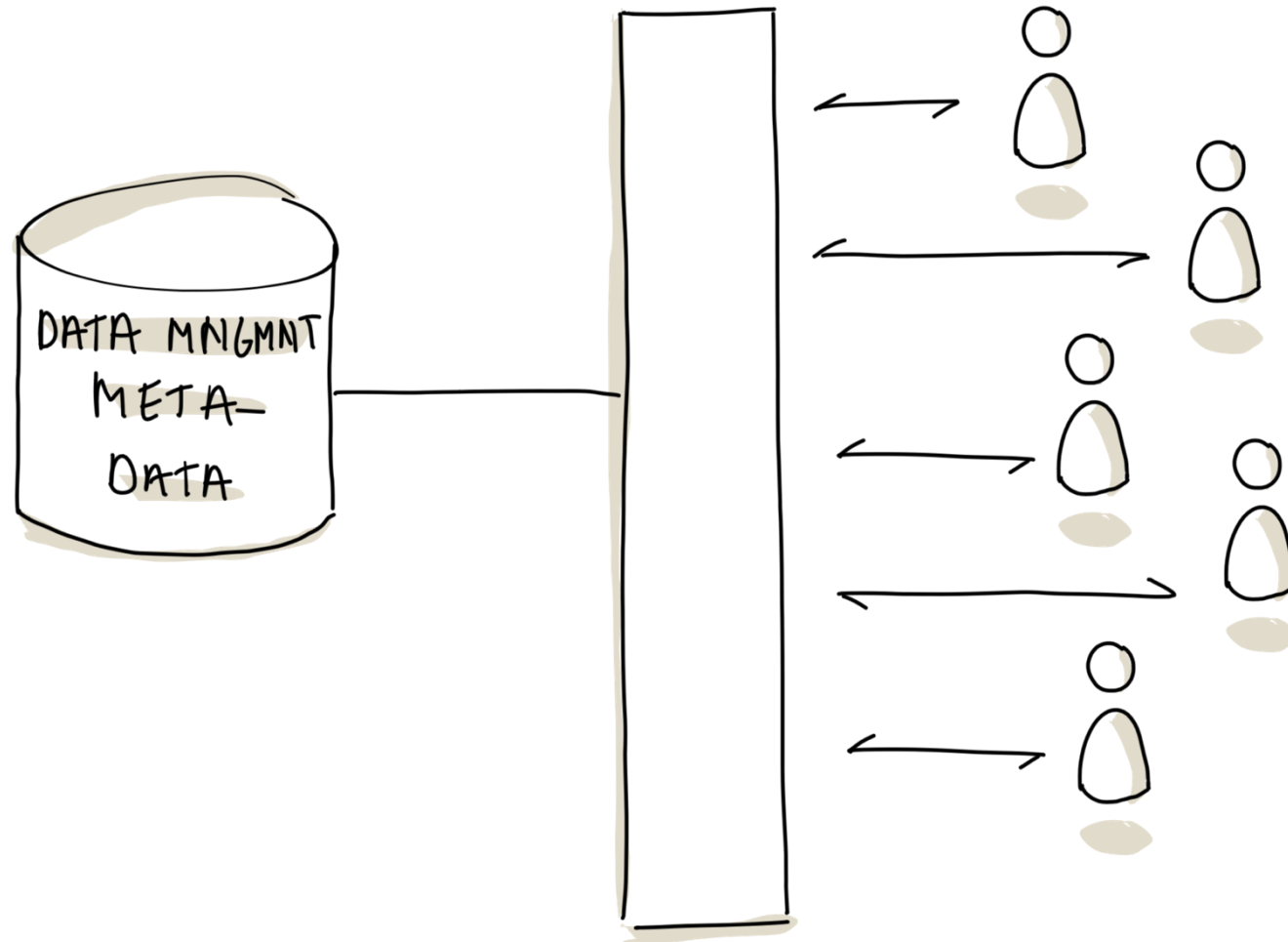
- Cohort analysis > Cohort analysis (Tableau View)**: A heatmap showing cohort analysis results.
- Cohort analysis > Cohort line graph (Tableau View)**: A line graph showing cohort analysis results.
- Cohort analysis > User transition example (Tableau View)**: A chart showing user transitions over time.
- Cohort analysis > User transition extra (Tableau View)**: A bar chart showing user transition data.
- Concept type dynamic changes > Concept types (Tableau View)**: A bar chart showing dynamics of concept changes.
- Concept type dynamic changes > Details (Tableau View)**: A table showing details of concept changes.
- DGC Analysis proxy logs > session by host (Tableau View)**: A bar chart showing session duration by host.
- DGC Analysis proxy logs > session duration (Tableau View)**: A horizontal bar chart showing session duration.
- DGC Analysis proxy logs > Sessions by country (Tableau View)**: A bar chart showing sessions by country.
- DGC Analysis proxy logs > sessions by host (Tableau View)**: A bar chart showing sessions by host.
- Foreign direct investment > Compare countries (Tableau View)**: A chart comparing foreign direct investment across countries.
- Foreign direct investment > Details (2000-2003) (Tableau View)**: A chart showing details of foreign direct investment from 2000 to 2003.

Each report card includes a 'Candidate' label and a green checkmark icon. The interface also features a 'Sort by Name' dropdown and a 'Validate' button. The page number '32' is visible in the bottom right corner.

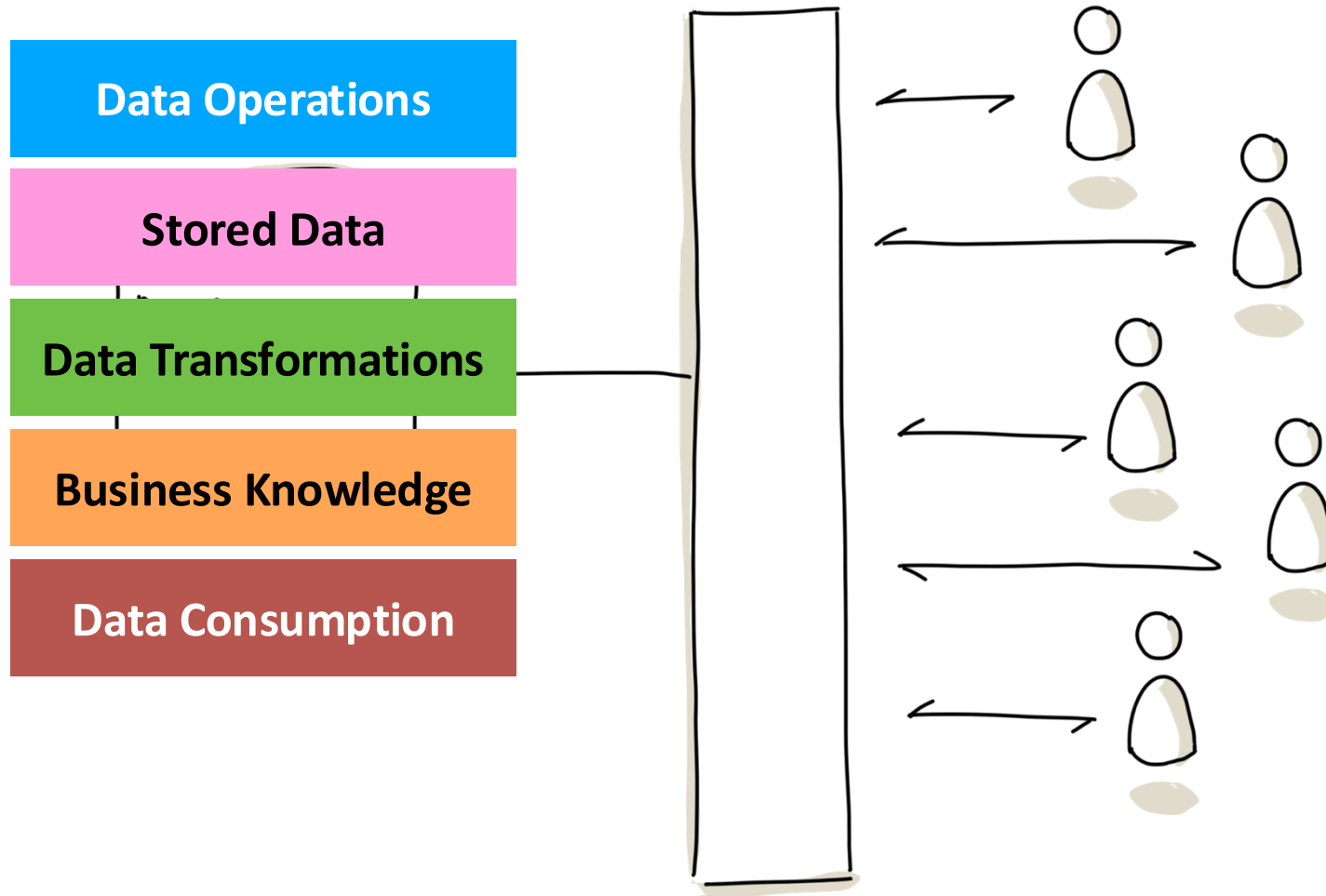
[Collibra (Screenshot deprecated – now integrated in the ‘Data Catalog’)]



META-DATA PORTAL



META-DATA PORTAL



META-DATA PORTAL

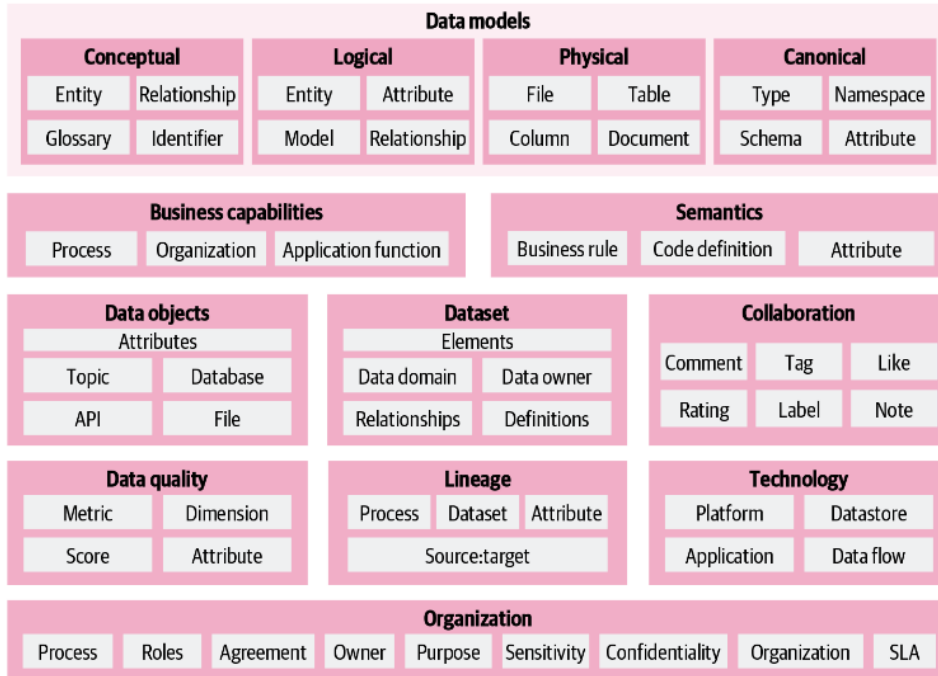
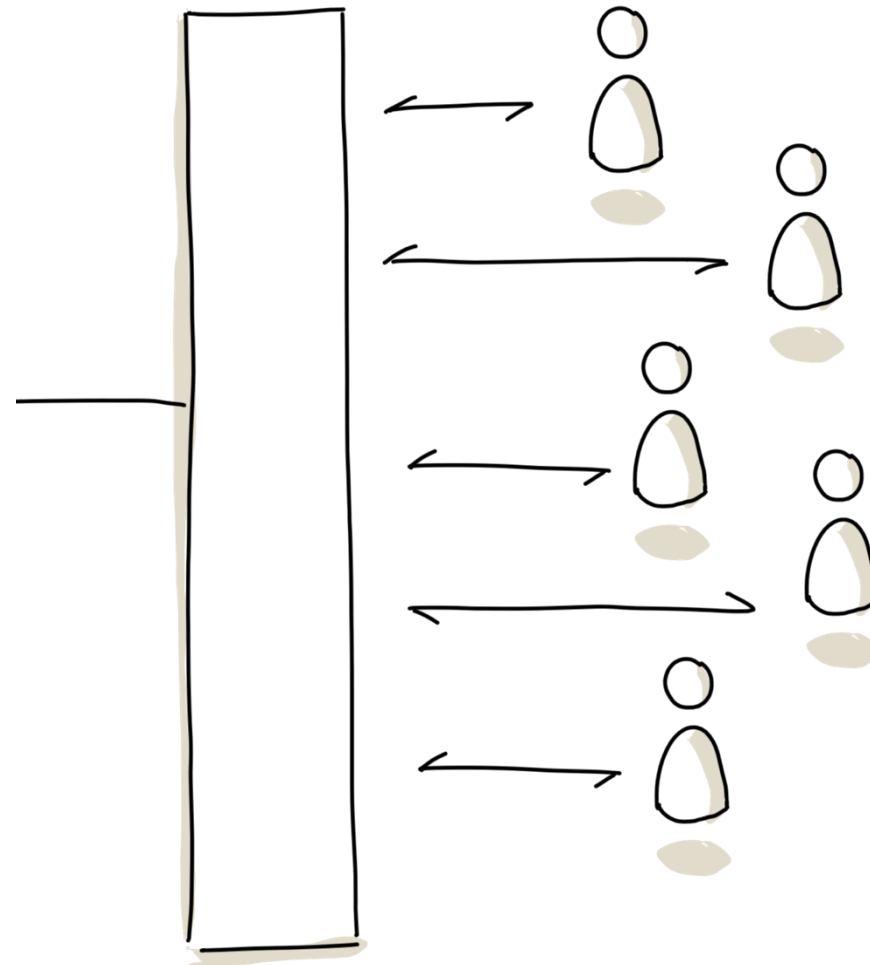


Figure 10-2. An overview showing all the major metadata management subject areas for the enterprise. These subject areas are nonexhaustive. Each organization uses the areas differently.



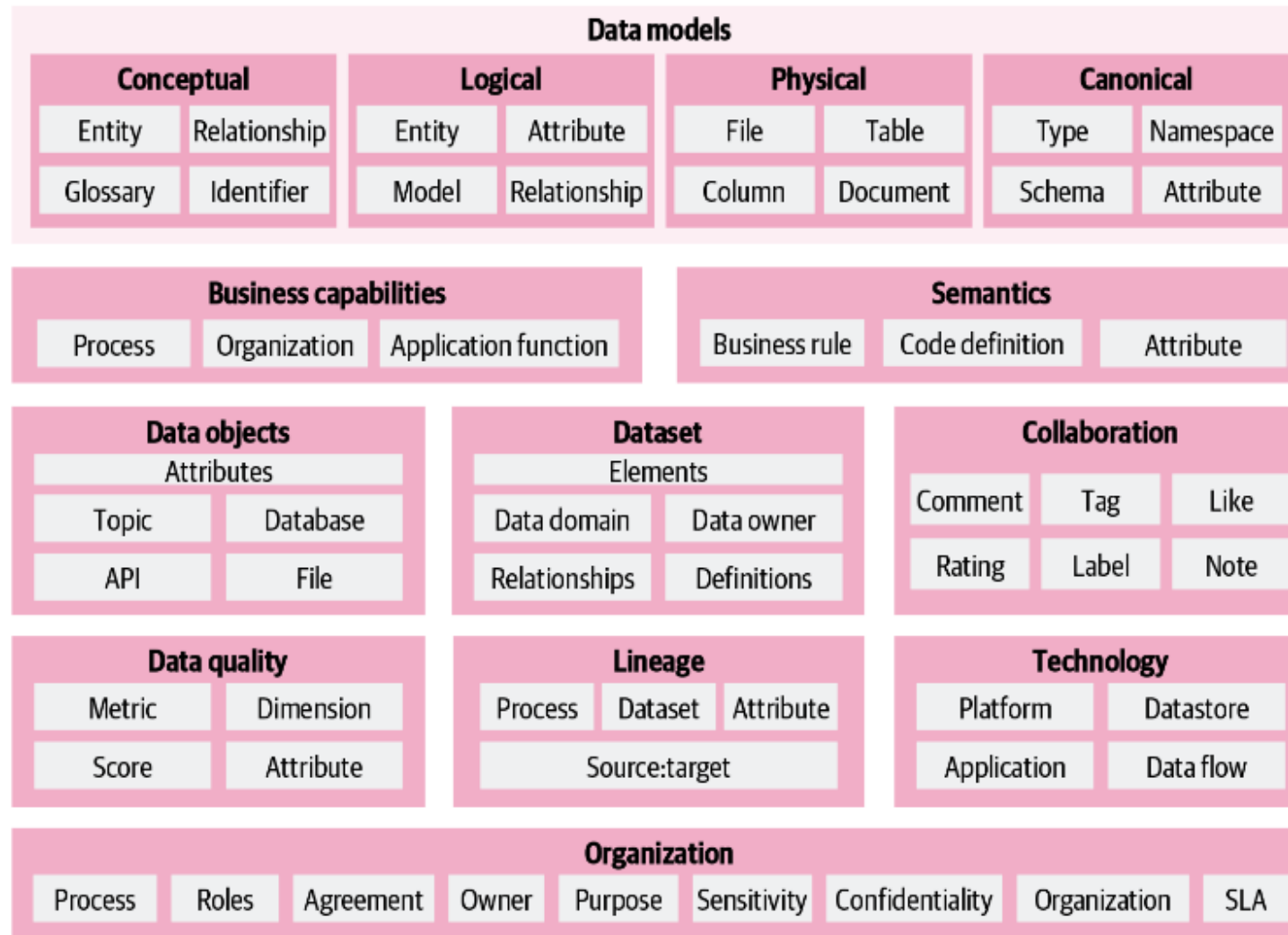
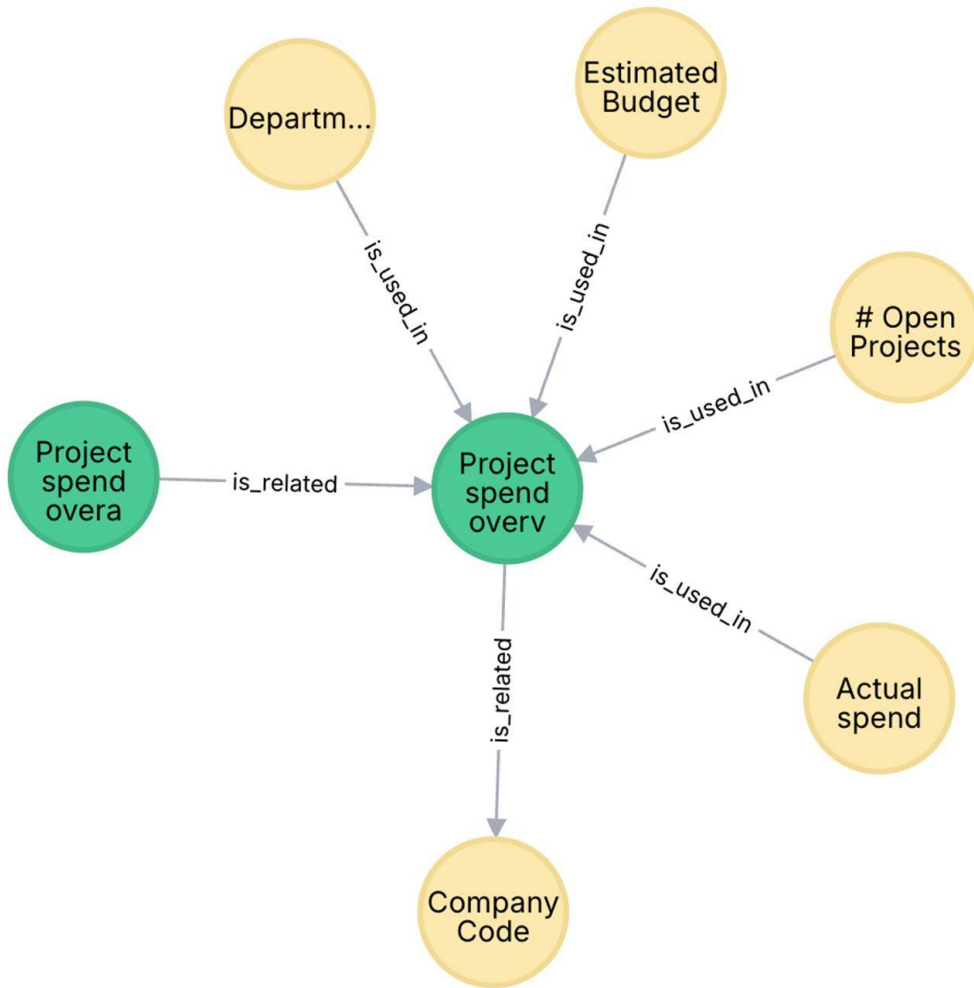


Figure 10-2. An overview showing all the major metadata management subject areas for the enterprise. These subject areas are nonexhaustive. Each organization uses the areas differently.





Asset Types

- Report
- Definition
- Webpage

Relation Type

- is used in / uses
- is related to
- is parent of / is child of
- mentions / is mentioned in



Start Discovering Your Data Assets

Search datasets, fields, visualizations, etc.



Topics 9

Collections of results to help you navigate through specific use cases.

Ap

Applications

Show all applications of our ecosystem which include a lineage to understand better...

BG

Business glossary

List all Glossary Items organized by Business Object and Business Data

DP

Data products

List all data products available in our organisation aligned with Data mesh approach

FD

Finance Domain

Lists all the related assets to the Finance Domain


Kp

KPIs

List all Key Performance Indicators available in our organisation

MD

Marketing domain

List all Items associated to Marketing domain 



Home > Knowledge Catalog

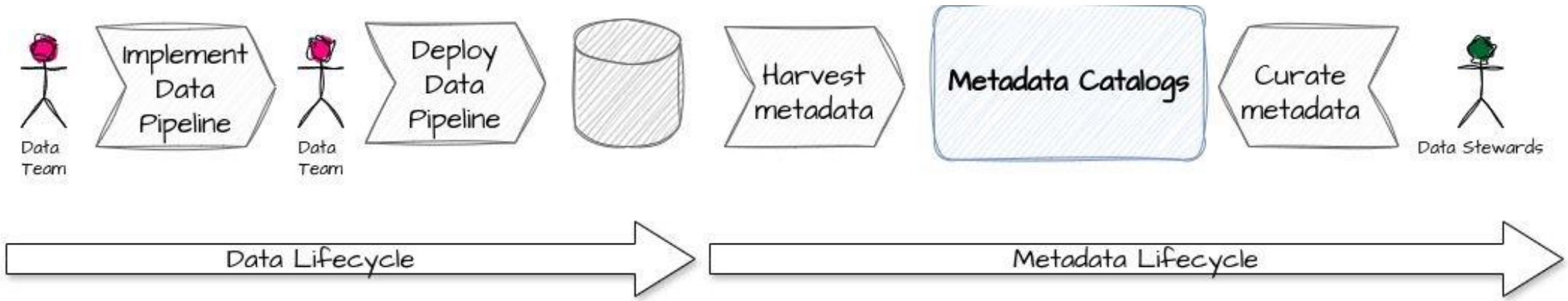
Data Assets

Filter by name, owner, creation date...

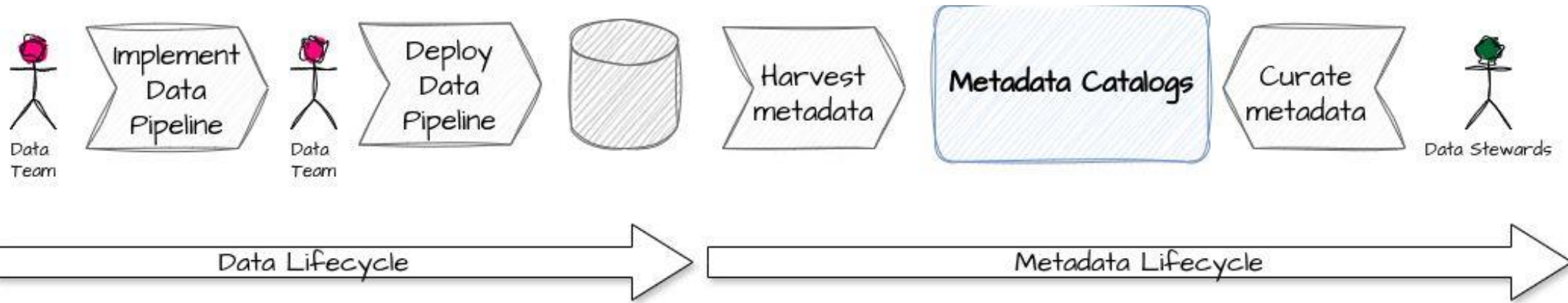
<input type="checkbox"/>	Name	Terms	Data Quality	# R
<input type="checkbox"/>	<u>src_person</u>	PII Employee Enum	<div><div style="width: 70%;"></div></div>	
<input type="checkbox"/>	<u>Master customer</u>	PII Customer	<div><div style="width: 30%;"></div></div>	
<input type="checkbox"/>	<u>Customers 2019</u>	PII Customer	<div><div style="width: 95%;"></div></div>	
<input type="checkbox"/>	<u>comp</u>	Account	<div><div style="width: 80%;"></div></div>	
<input type="checkbox"/>	<u>Customer campaigns</u>	Customer Campaign	<div><div style="width: 95%;"></div></div>	
<input type="checkbox"/>	<u>cstmr</u>	PII Customer	<div><div style="width: 95%;"></div></div>	
<input type="checkbox"/>	<u>employees_2020</u>	PII Employee	<div><div style="width: 70%;"></div></div>	
<input type="checkbox"/>	<u>Master address</u>	Address	<div><div style="width: 30%;"></div></div>	
<input type="checkbox"/>	<u>cstomers_2019_ext</u>	PII Customer	<div><div style="width: 95%;"></div></div>	
<input type="checkbox"/>	<u>account_list</u>	PII Account	<div><div style="width: 80%;"></div></div>	



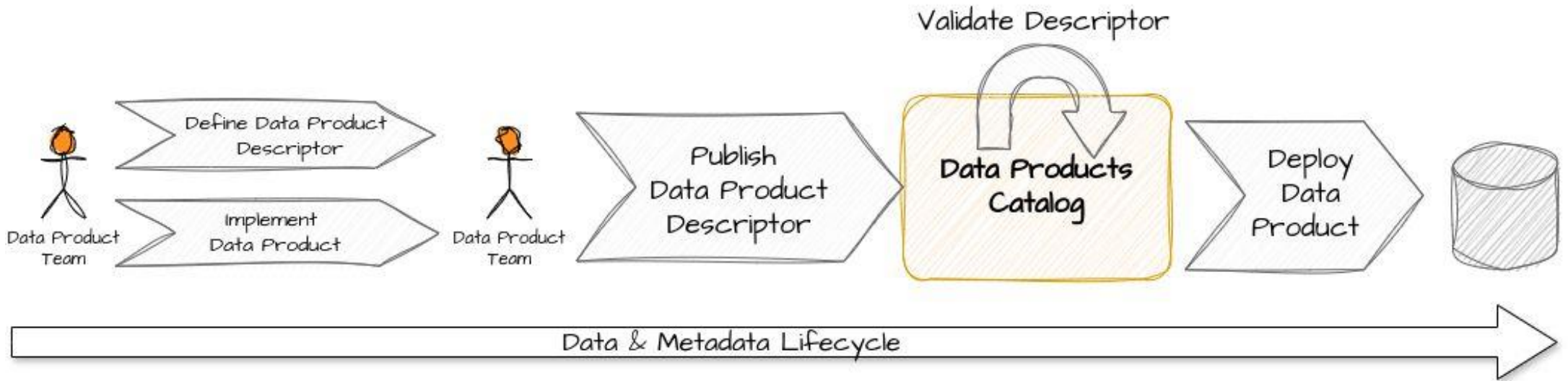
CENTRALIZED



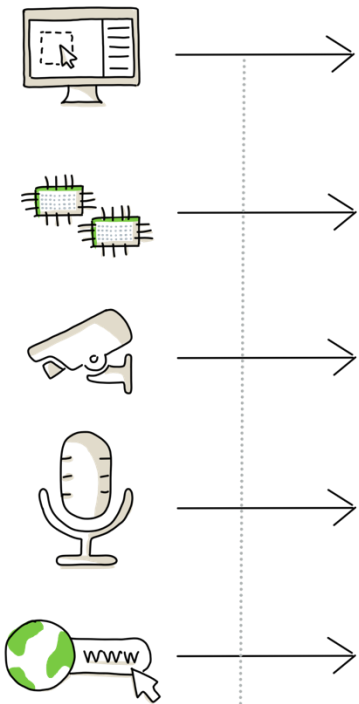
CENTRALIZED



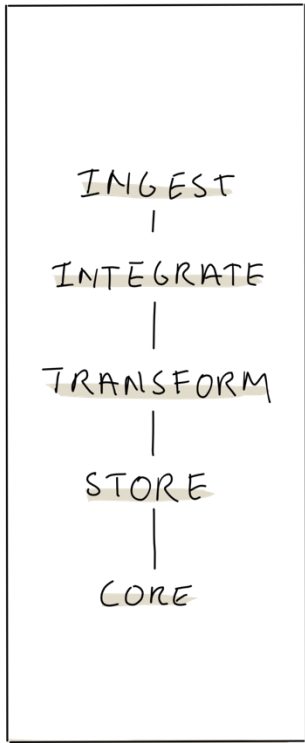
DATA MESH



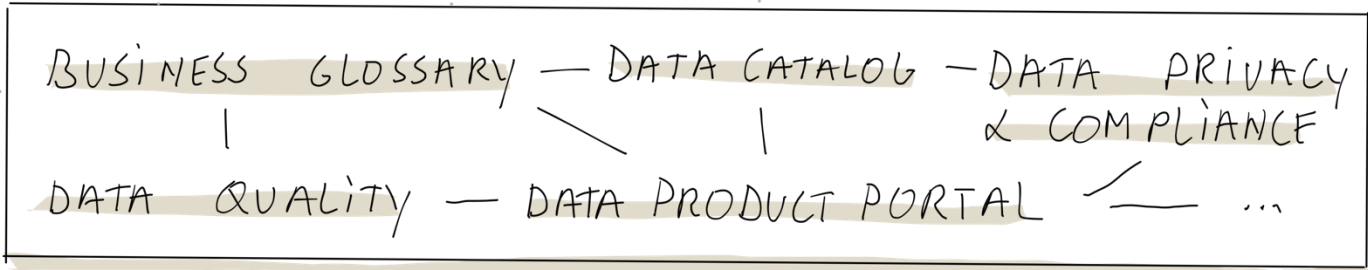
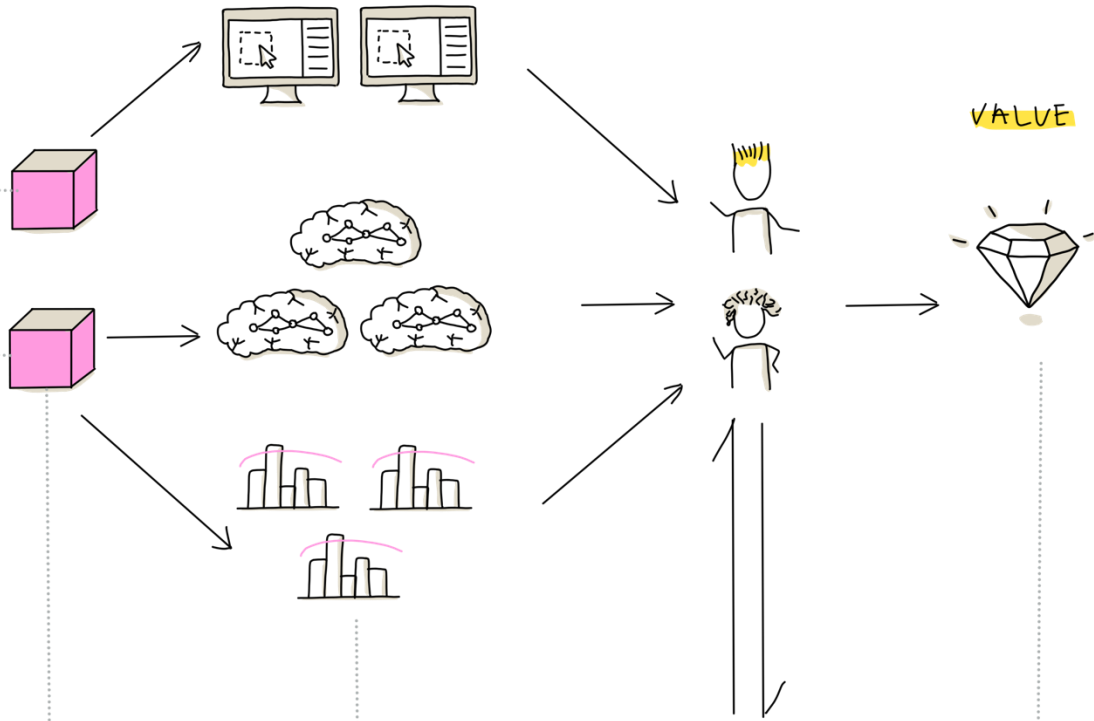
DATA PRODUCERS (SOURCES)



DATA PLATFORM



CONSUMPTION



META-DATA MANAGEMENT



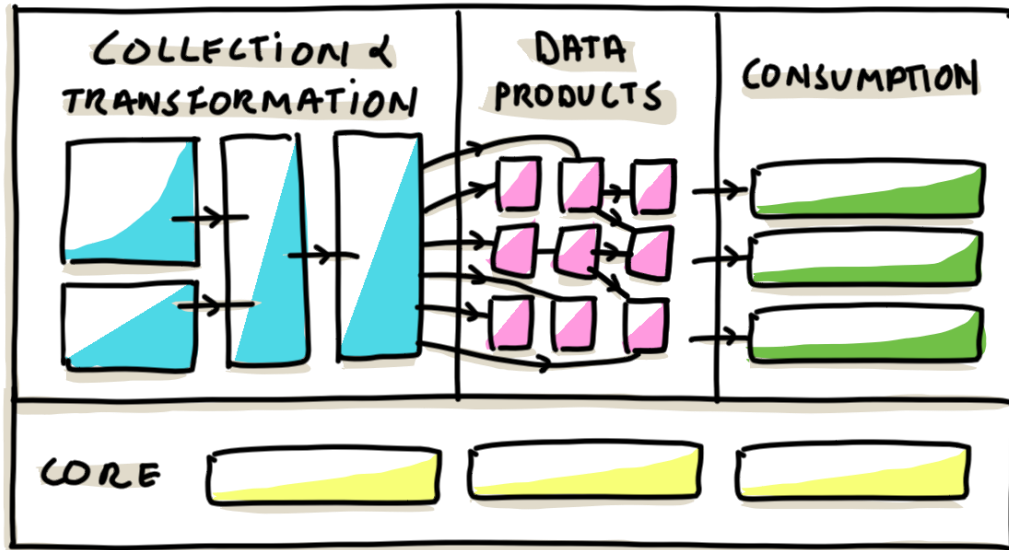
Table of Contents

- Dead Horse Theory
- Data Platform
 - Introduction
 - Core Layers
 - Additional Layers
- **Technology Selection**

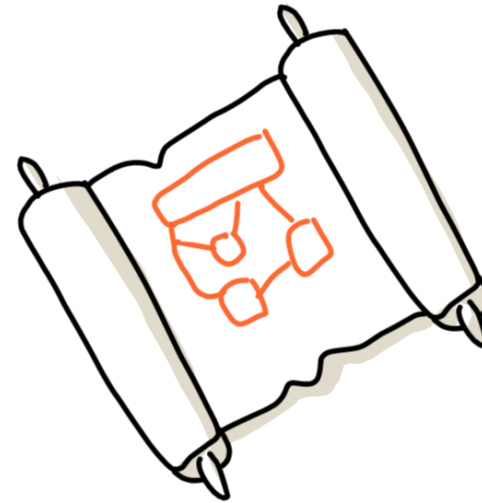


From Platform to Implementation

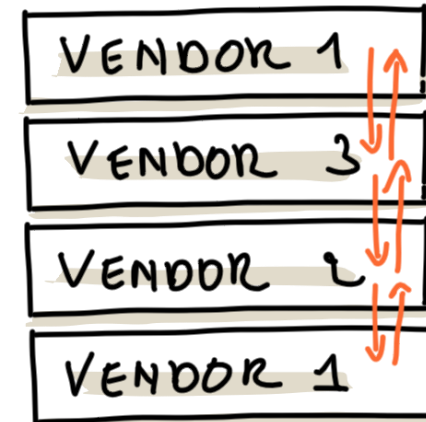
DATA PLATFORM



ARCHITECTURE



DATA STACK



Medaillon Reference Architecture



LAKEHOUSE STORAGE LAYERS

SOURCES



BRONZE



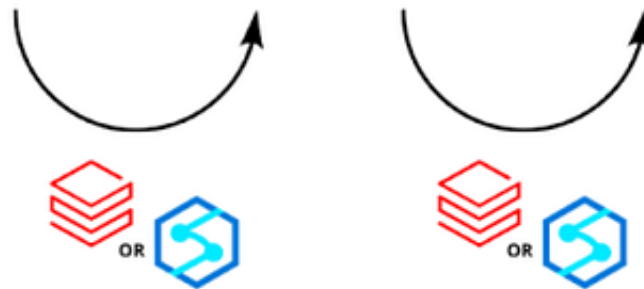
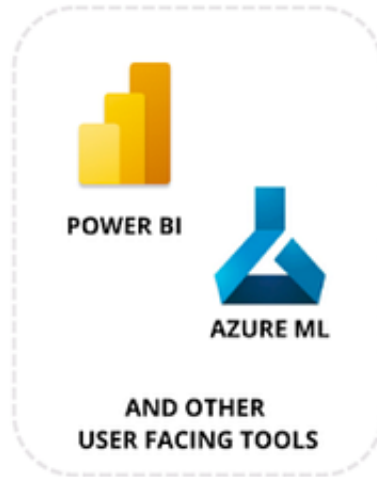
SILVER



GOLD

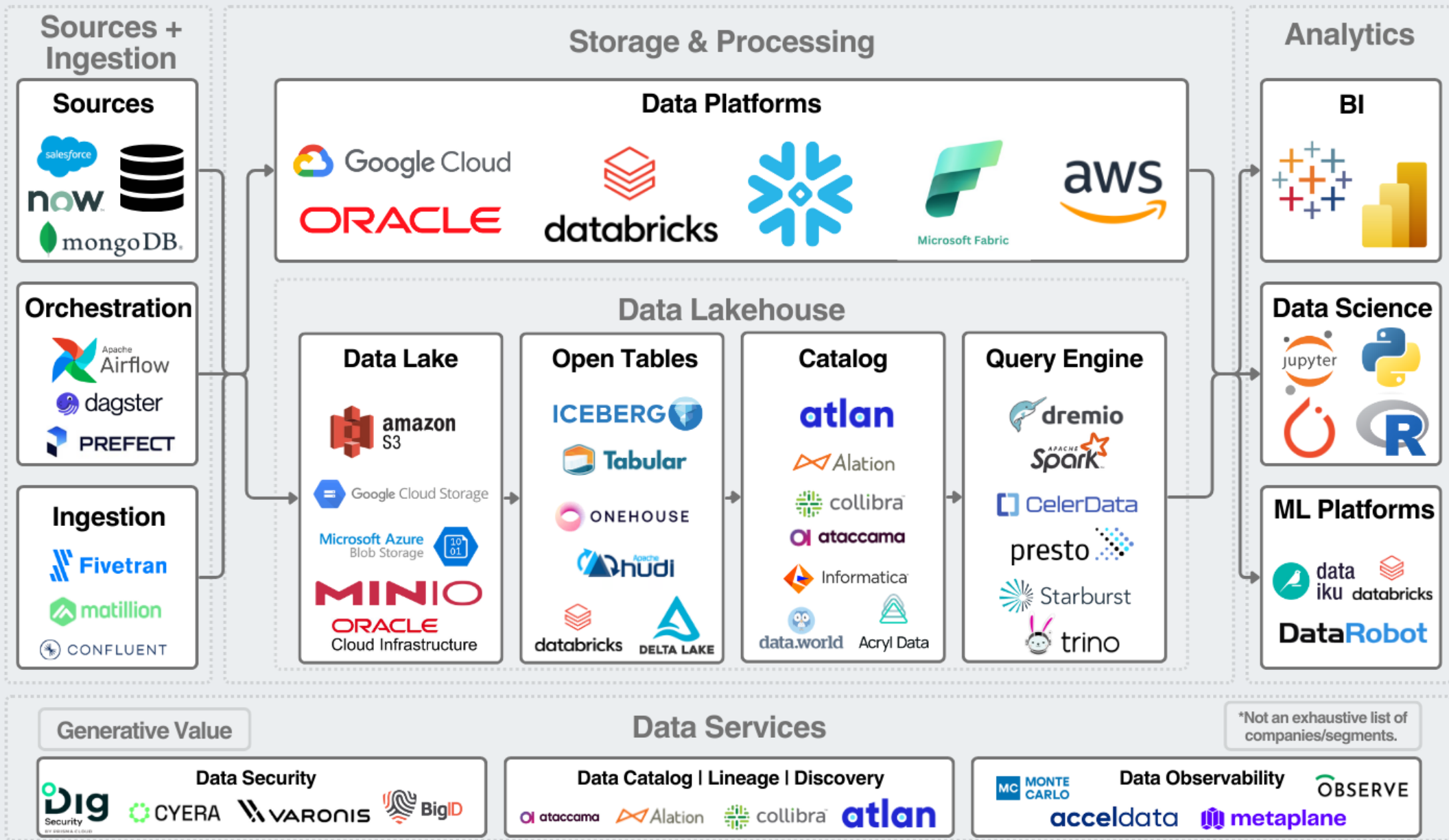


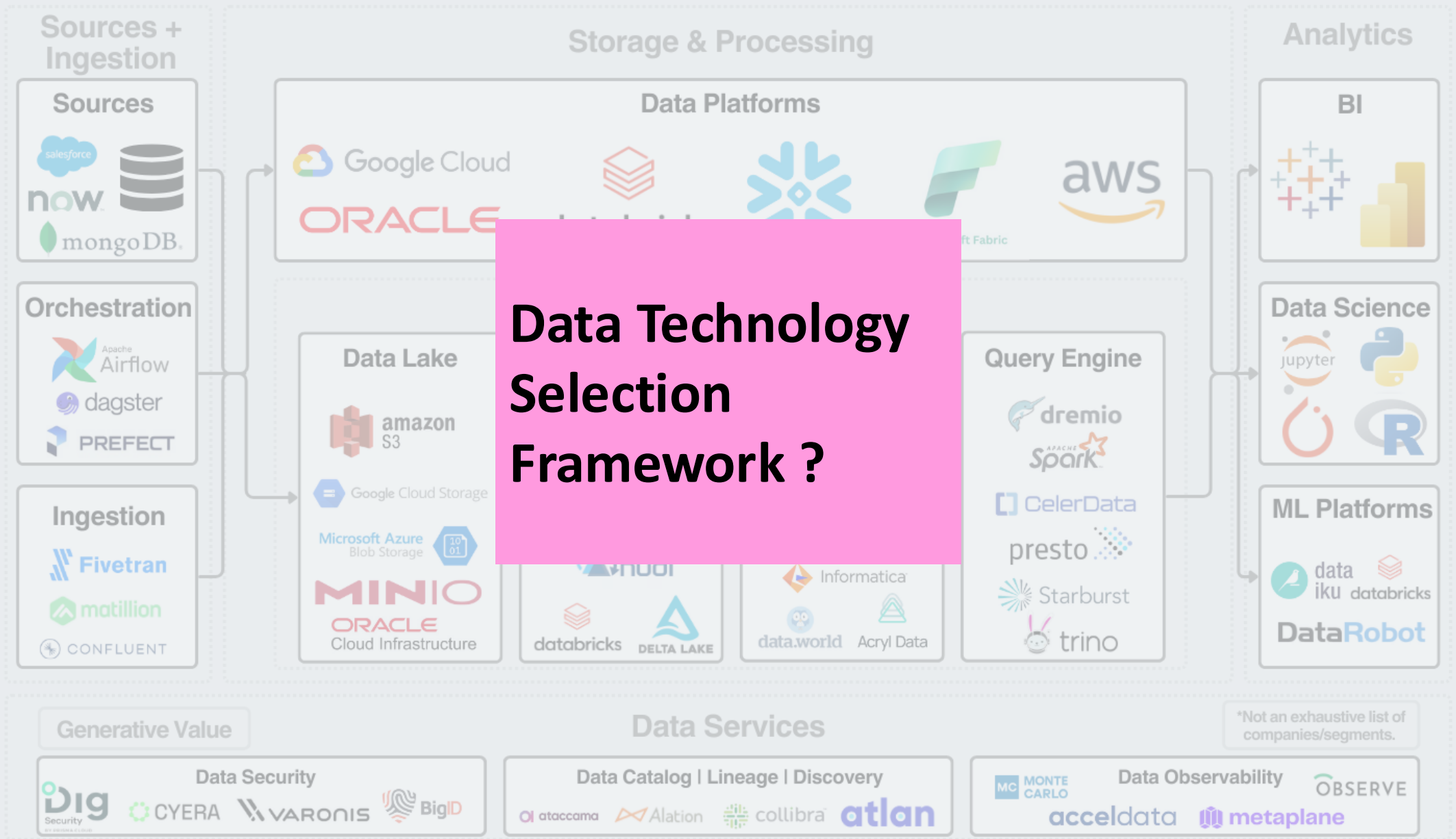
SERVE



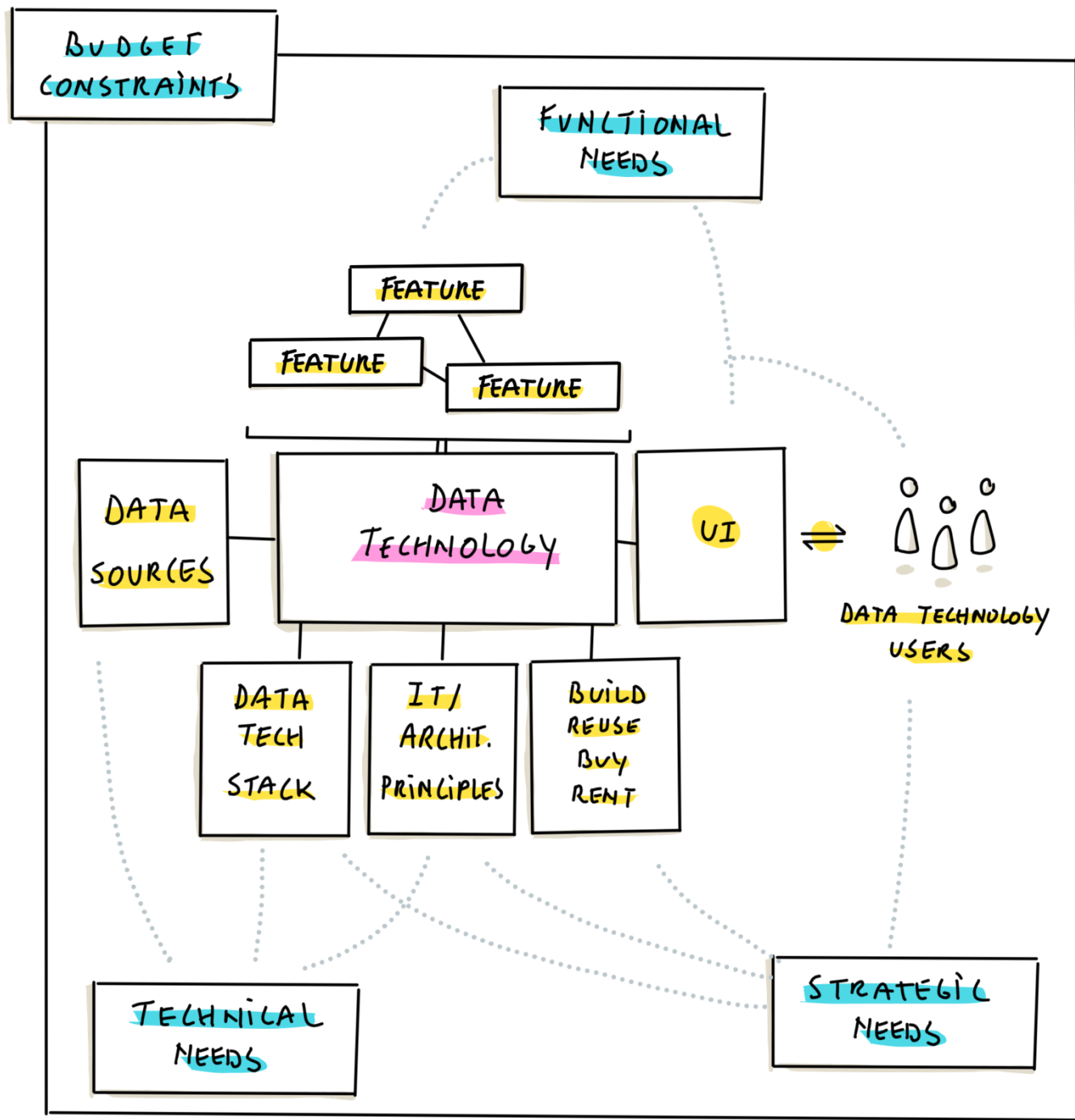
PROCESSING AND TRANSFORMATION WITH DATABRICKS OR SYNAPSE PIPELINES







Data Technology Selection Framework



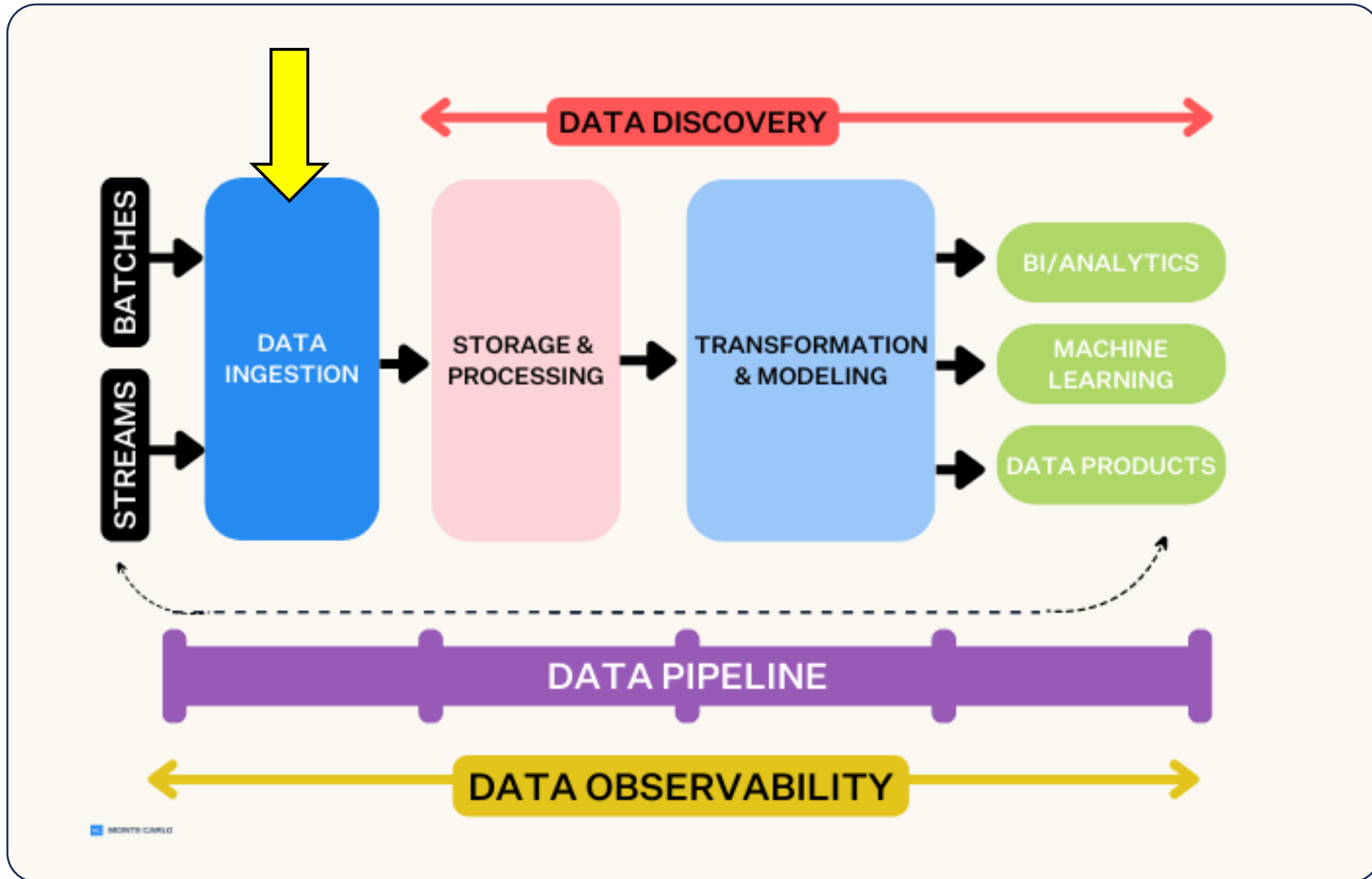
**Data Technology
Selection
Framework**

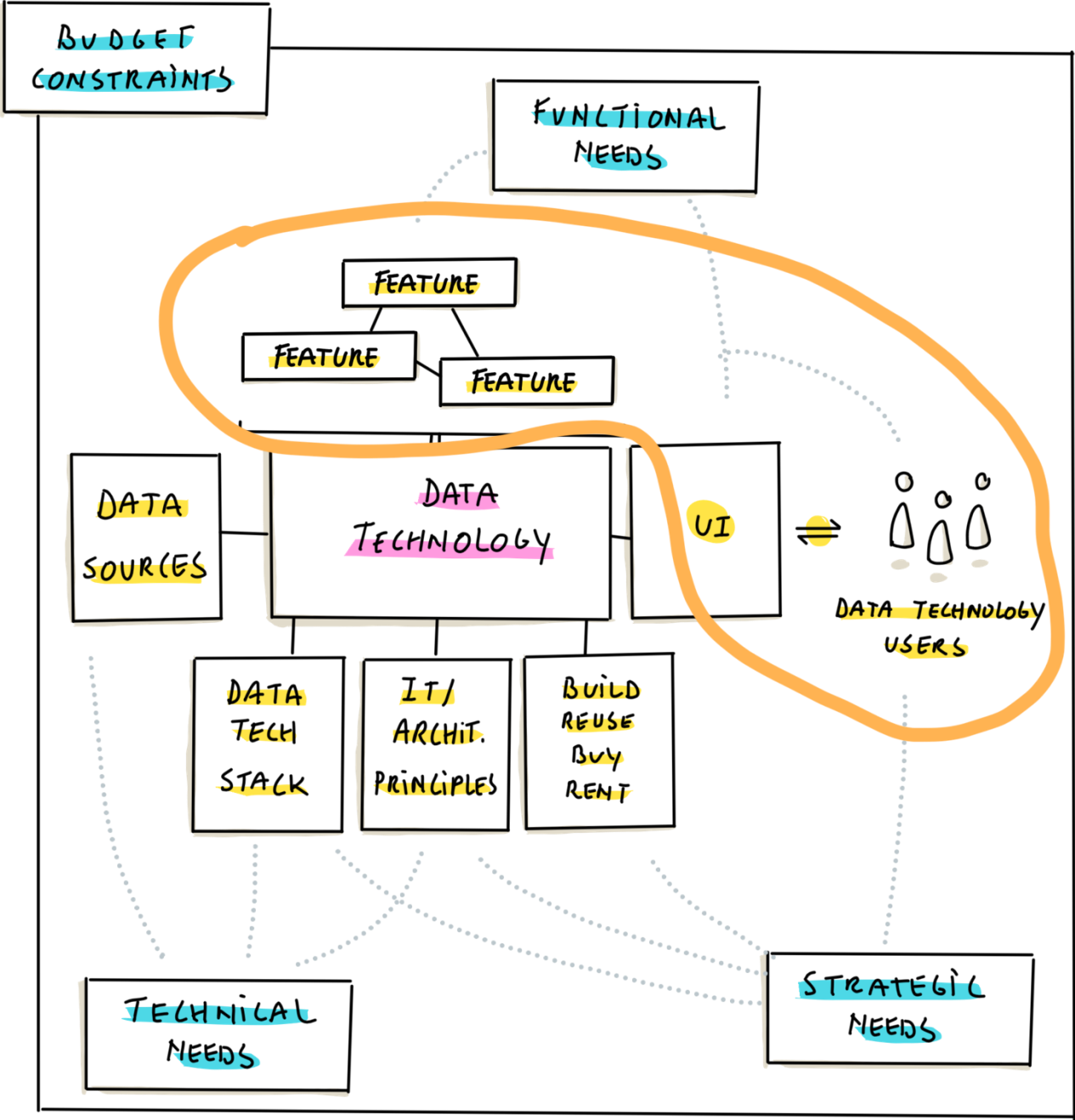
**Complete Data
Platform**

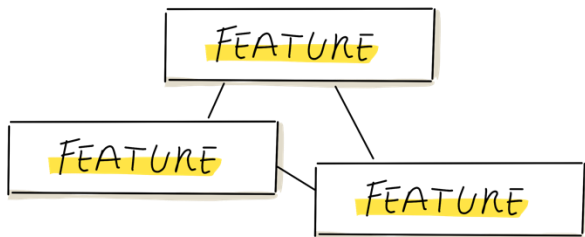
**1 or More Data
Platform
Component(s)**



Example: Data Ingestion



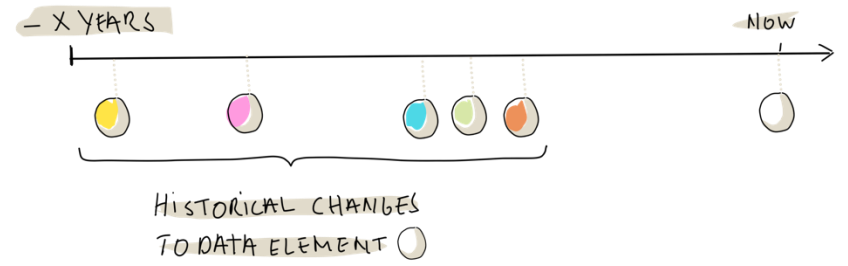




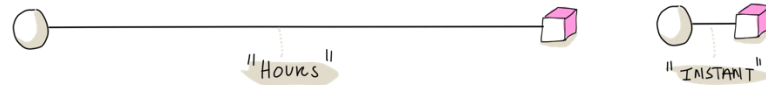
DATA FRESHNESS



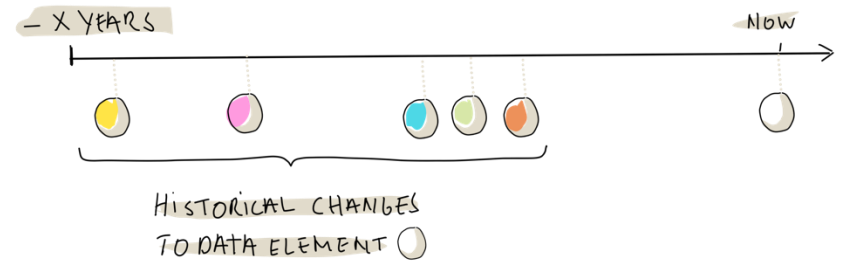
TIME TRAVELING



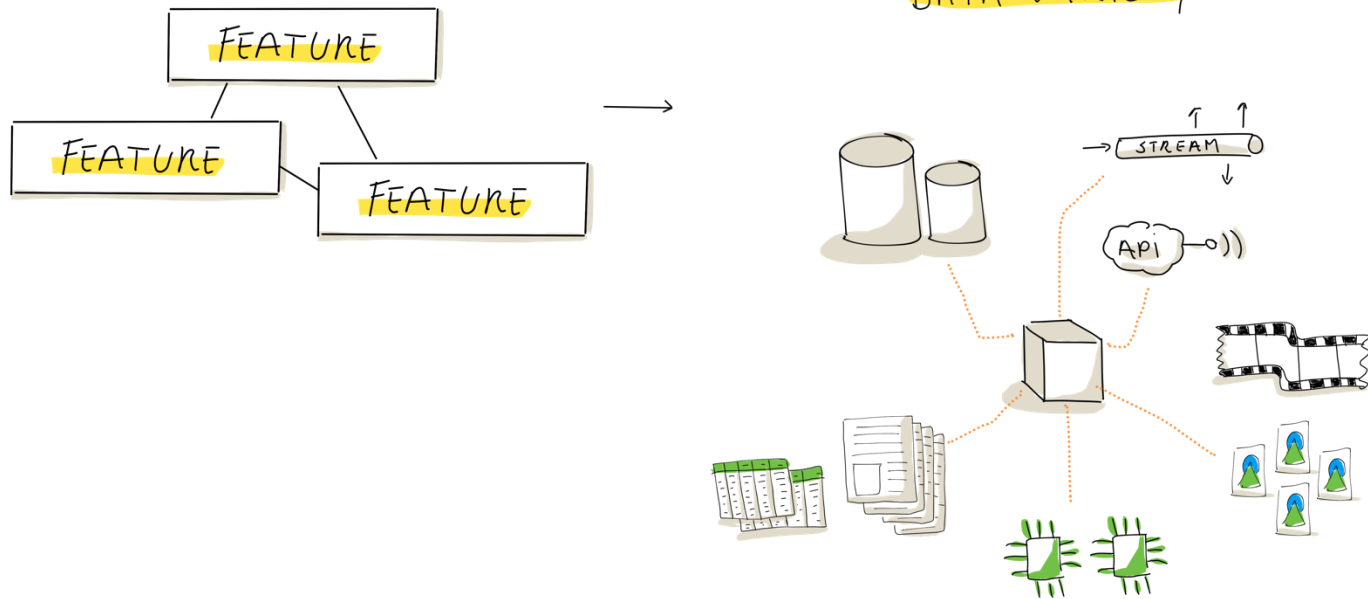
DATA FRESHNESS



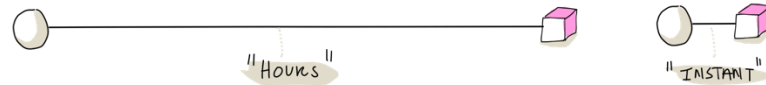
TIME TRAVELING



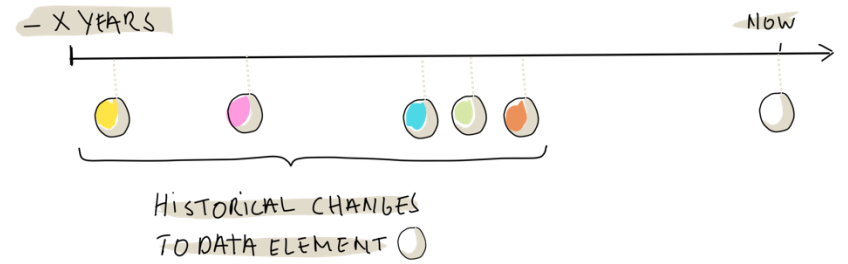
DATA VARIETY



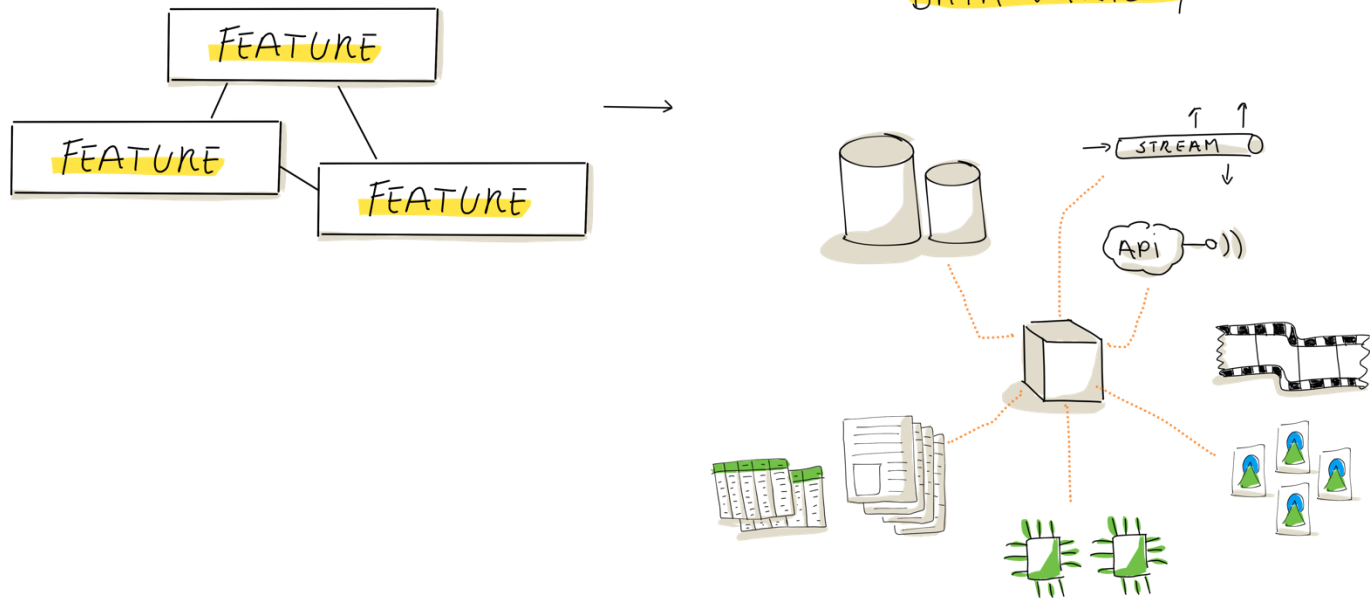
DATA FRESHNESS



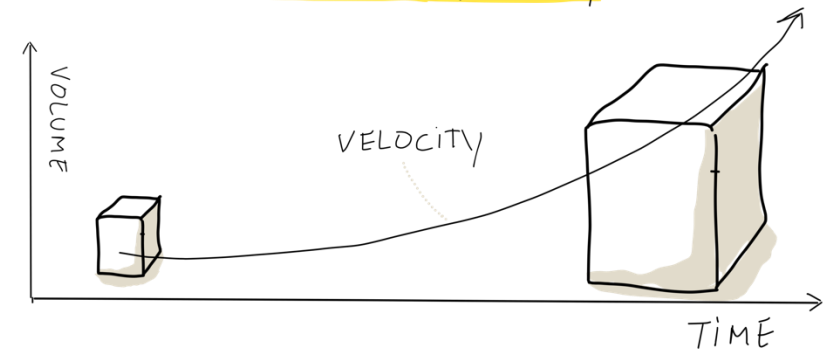
TIME TRAVELING



DATA VARIETY



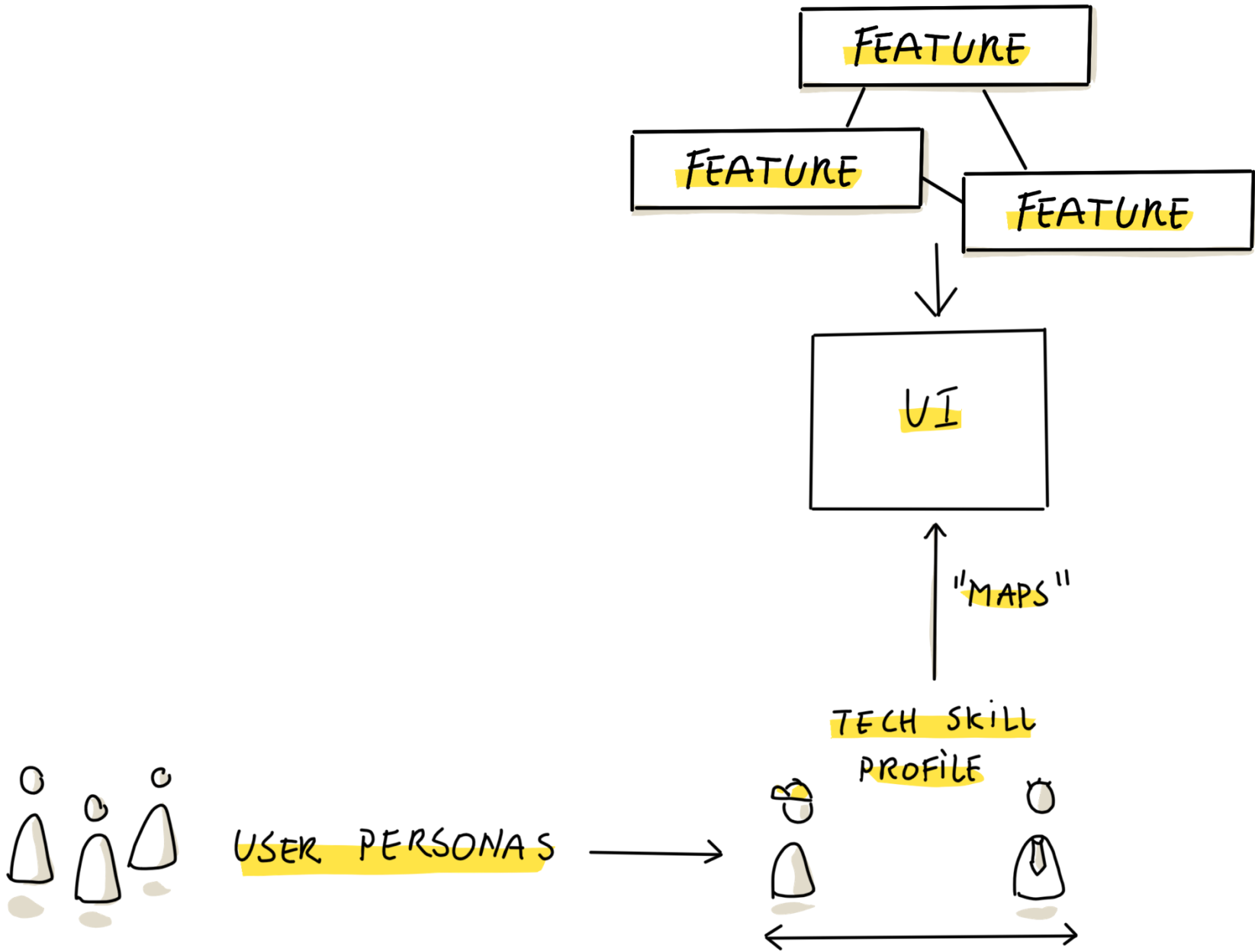
DATA VOLUME & VELOCITY

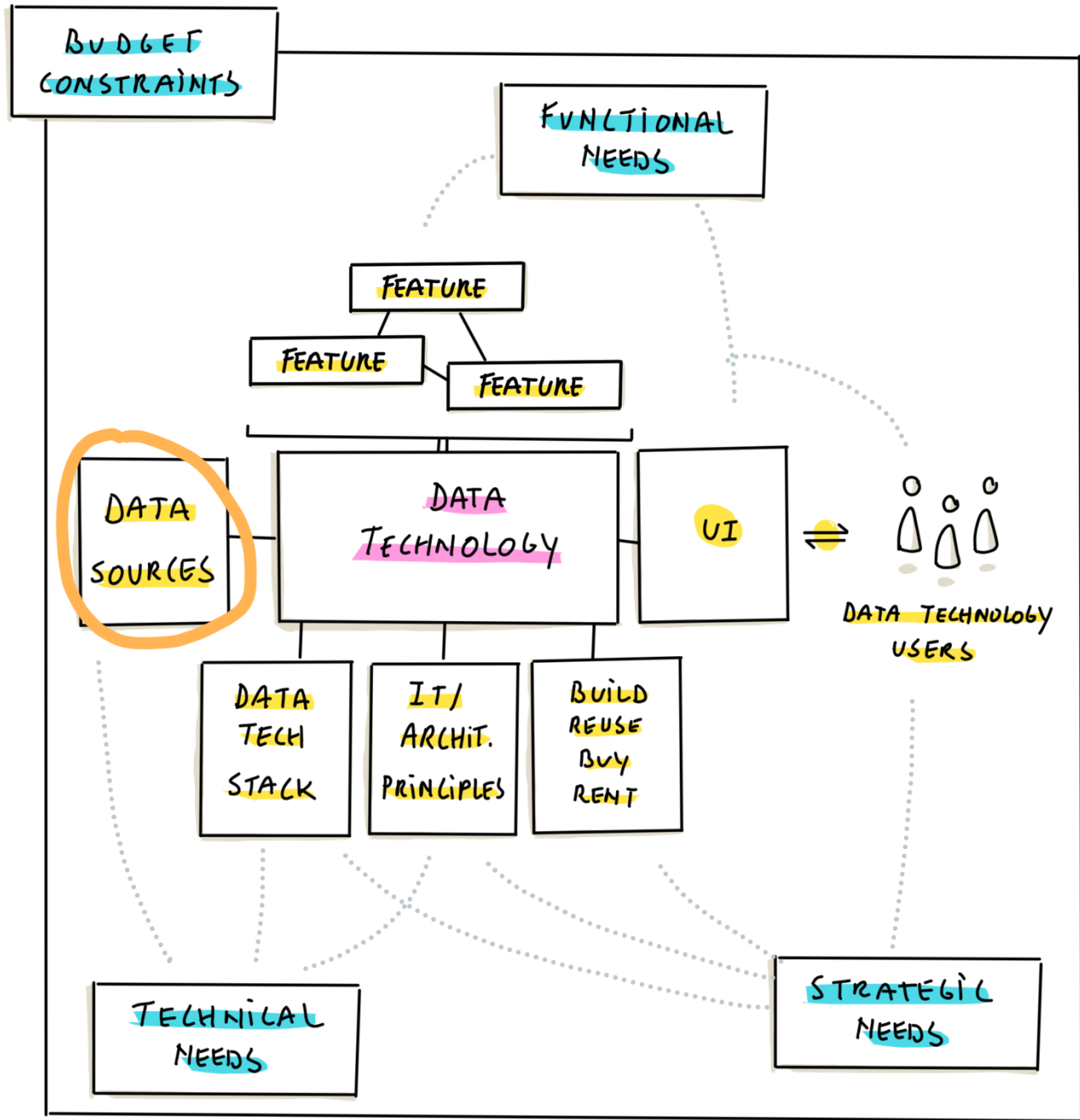


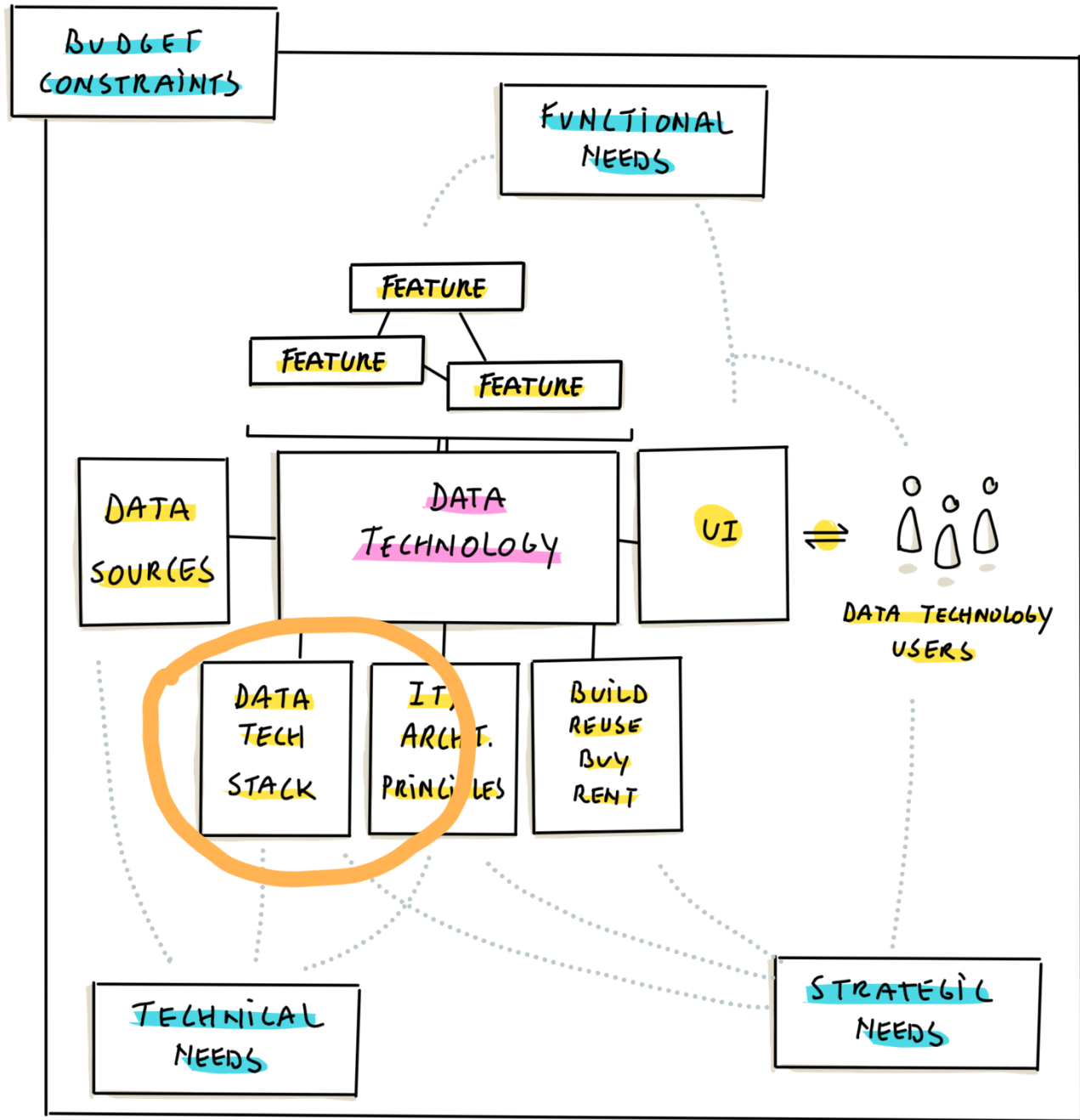


USER PERSONAS



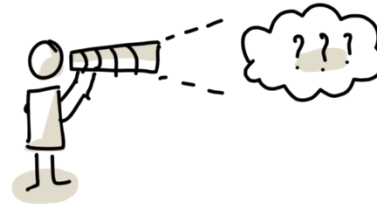








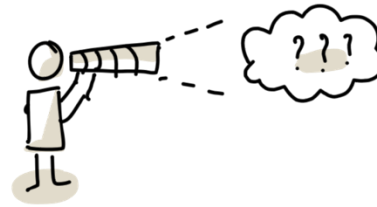
DATA PLATFORM
ARCHITECTURE



IT
STRATEGY



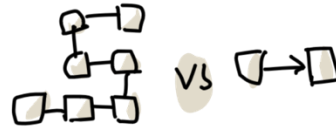
DATA PLATFORM
ARCHITECTURE



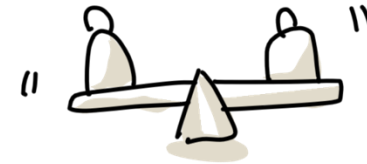
IT
STRATEGY



VENDOR
LOCK-IN



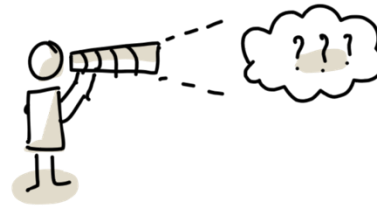
SIMPLICITY



MAINTENANCE
VS
CONTROL



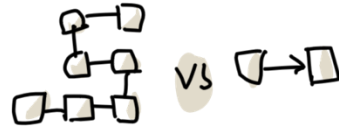
DATA PLATFORM
ARCHITECTURE



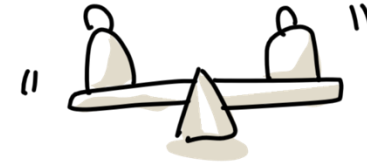
IT
STRATEGY



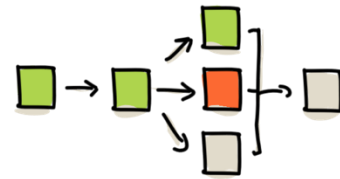
VENDOR
LOCK-IN



SIMPLICITY



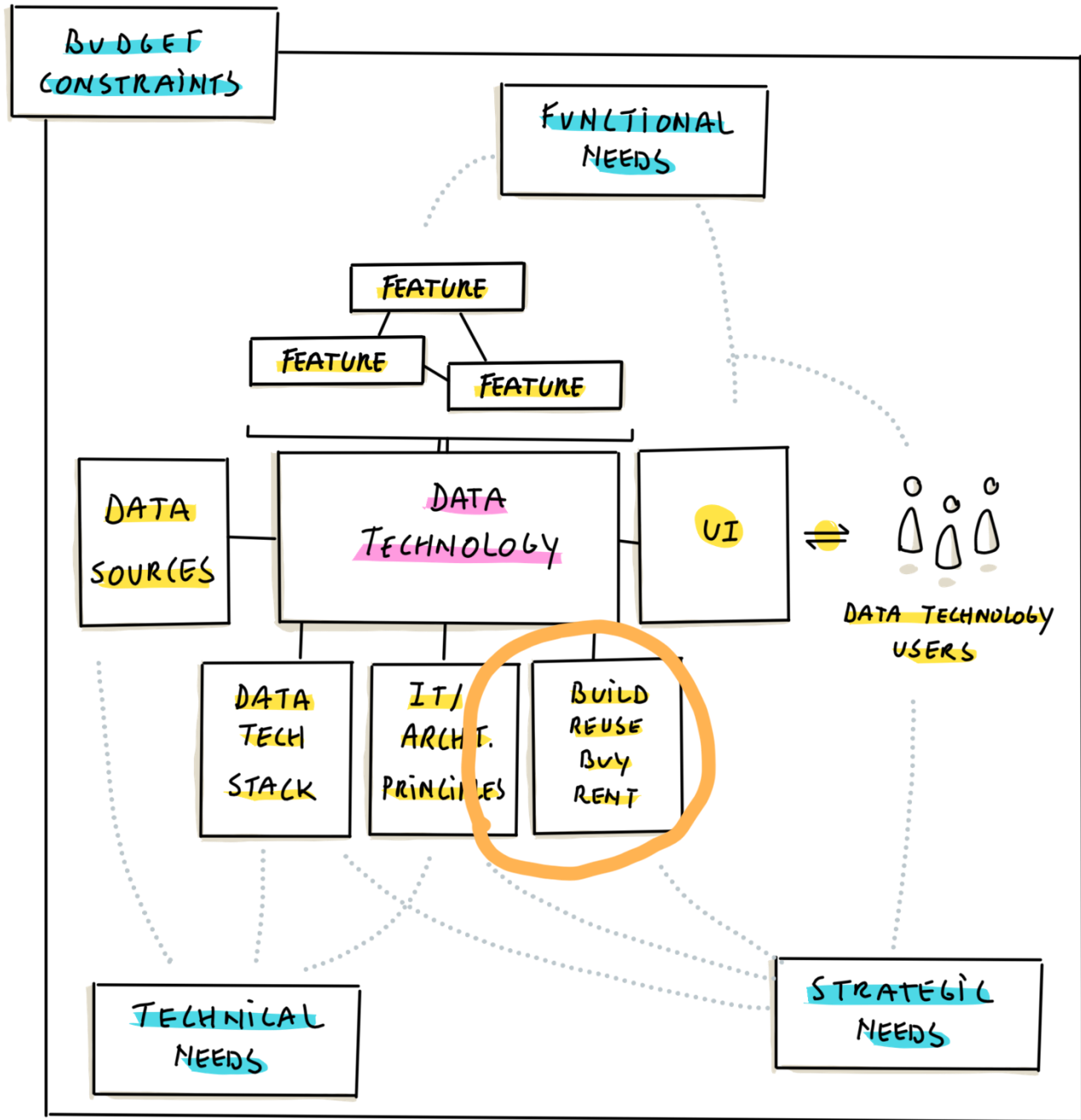
MAINTENANCE
VS
CONTROL

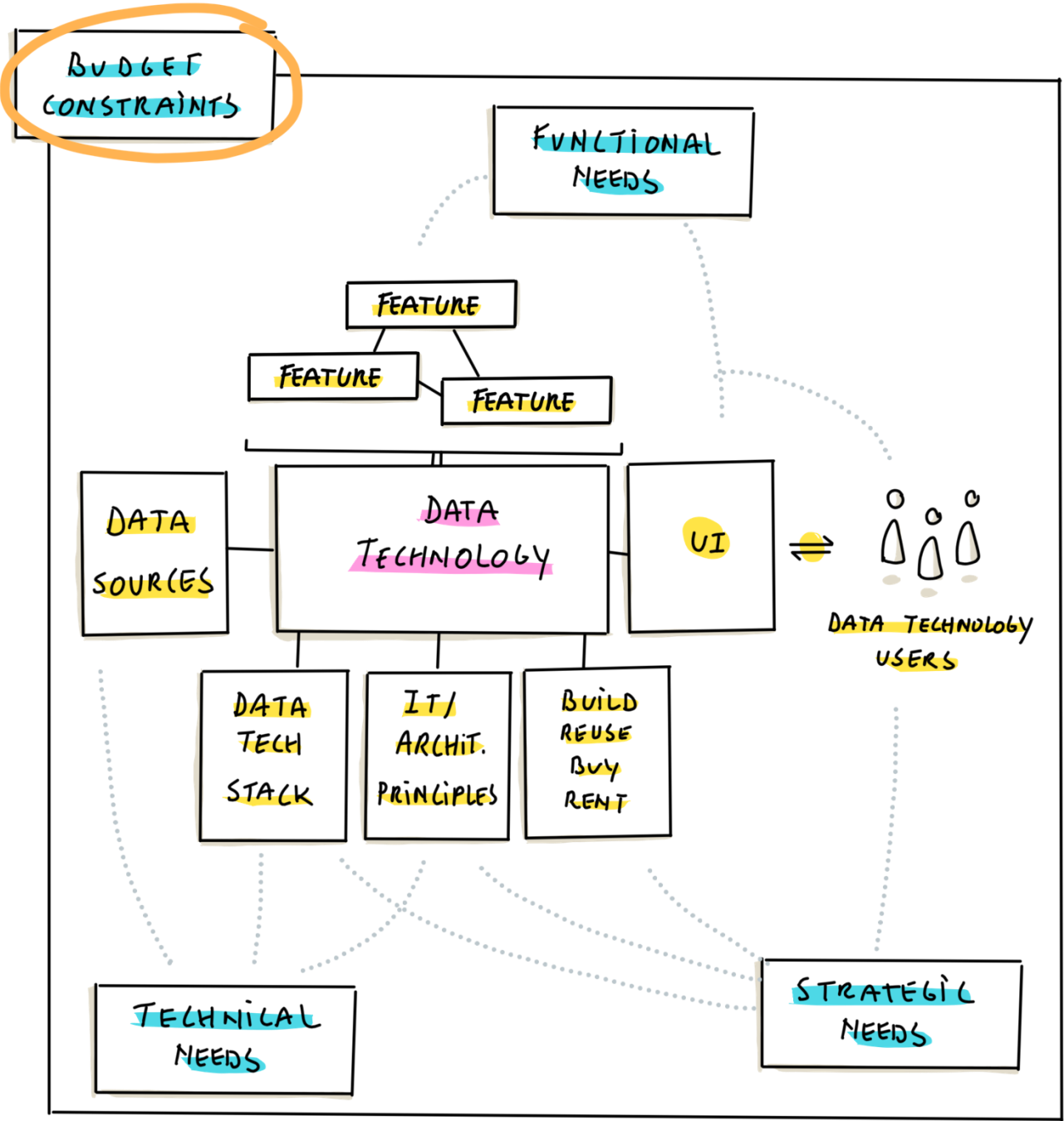


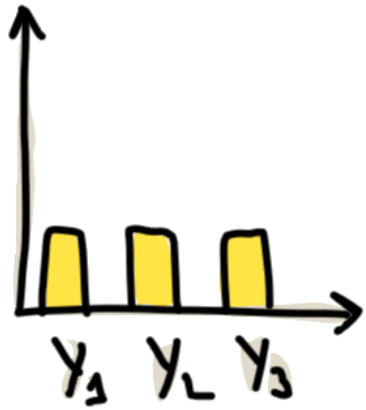
ORCHESTRATION
& MONITORING



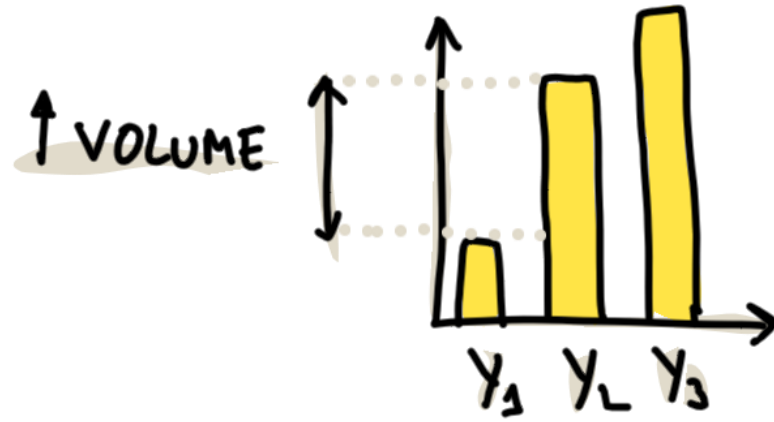
TOOL
ECOSYSTEM



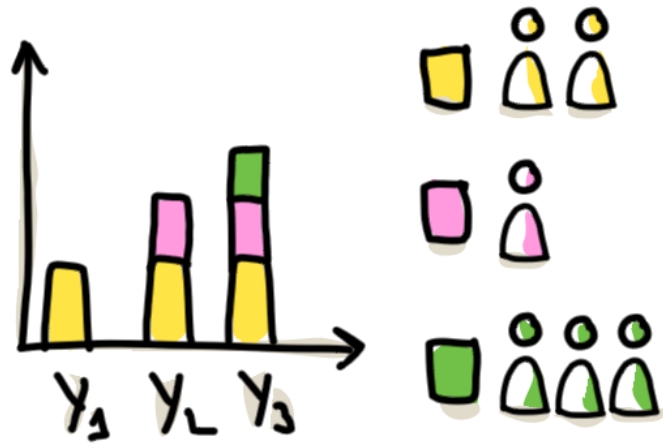




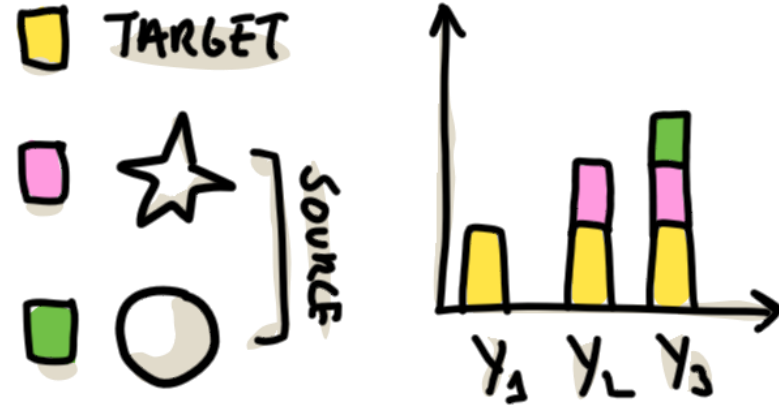
FIXED PRICE



PER VOLUME



PER USER



PER SOURCE/TARGET

PART	CRITERIUM	DESCRIPTION
Functional / Strategic	Future users	What will be the main profile of the future technology users?
Functional	Features	Which technology features are important?
Technical	Data Sources	What are the exact properties of the data sources?
Technical	Data Tech Stack	In which data architecture should the new technology be implemented?
Technical	Eco System	Which data tool eco-system is in use and will we keep using?
Technical / Strategic	IT Strategy	Cloud? Open Source? Micro Services? ...
Technical / Strategic	Simplicity	Choos the most simple (= direct) solution
Technical / Strategic	Vendor Lock-In	Avoid lock-in
Technical / Strategic	Monitoring / Orchestration	Integration in overall monitoring / orchestration
Budget Constraints	Pricing strategy	
Budget Constraints	TCO	



PART	CRITERIUM	DESCRIPTION
Purchasing / Strategic	Turnover vendor	What's the turnover of the vendor in the last X years?
Purchasing / Strategic	Support in EU	Is there a support organization in the EU?
Purchasing / Strategic	Product Maturity & Stability	What's the maturity of the product? (Can also be evaluated for OS)
Purchasing / Strategic	Community	What's that response time to resolve critical issues of an OS product?



EXERCISE : SELECTION CRITERIA WITHIN YOUR USE CASE?

Pick a needed
technology



PART	CRITERIUM	DESCRIPTION
Functional / Strategic	Future users	What will be the main profile of the future technology users?
Functional	Features	Which technology features are important?
Technical	Data Sources	What are the exact properties of the data sources?
Technical	Data Tech Stack	In which data architecture should the new technology be implemented?
Technical	Eco System	Which data tool eco-system is in use and will we keep using?
Technical / Strategic	IT Strategy	Cloud? Open Source? Micro Services? ...
Technical / Strategic	Simplicity	Choos the most simple (= direct) solution
Technical / Strategic	Vendor Lock-In	Avoid lock-in
Technical / Strategic	Monitoring / Orchestration	Integration in overall monitoring / orchestration
Budget Constraints	Pricing strategy	
Budget Constraints	TCO	



How to use in Practice

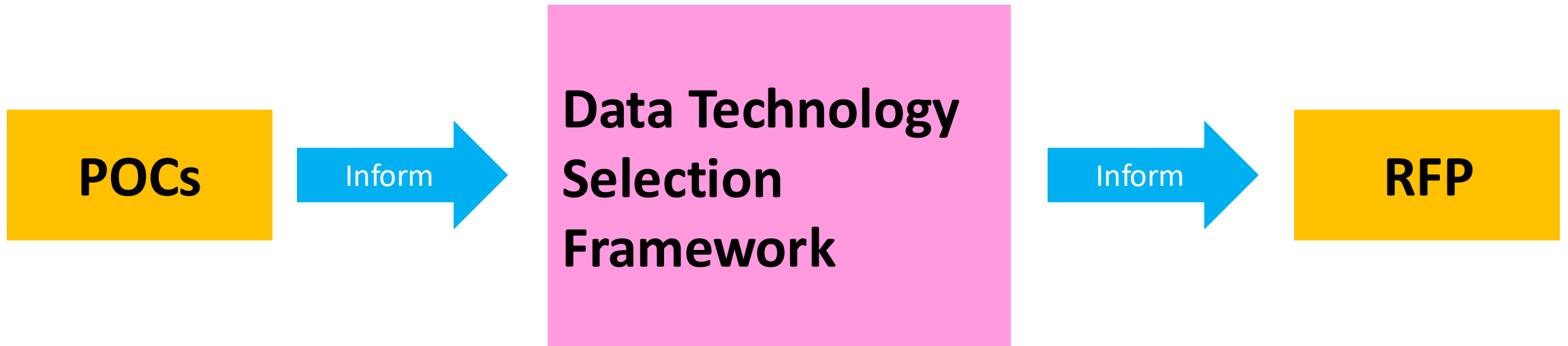


Table of Contents

- Dead Horse Theory
- Data Platform
 - Introduction
 - Core Layers
 - Additional Layers
- Technology Selection

