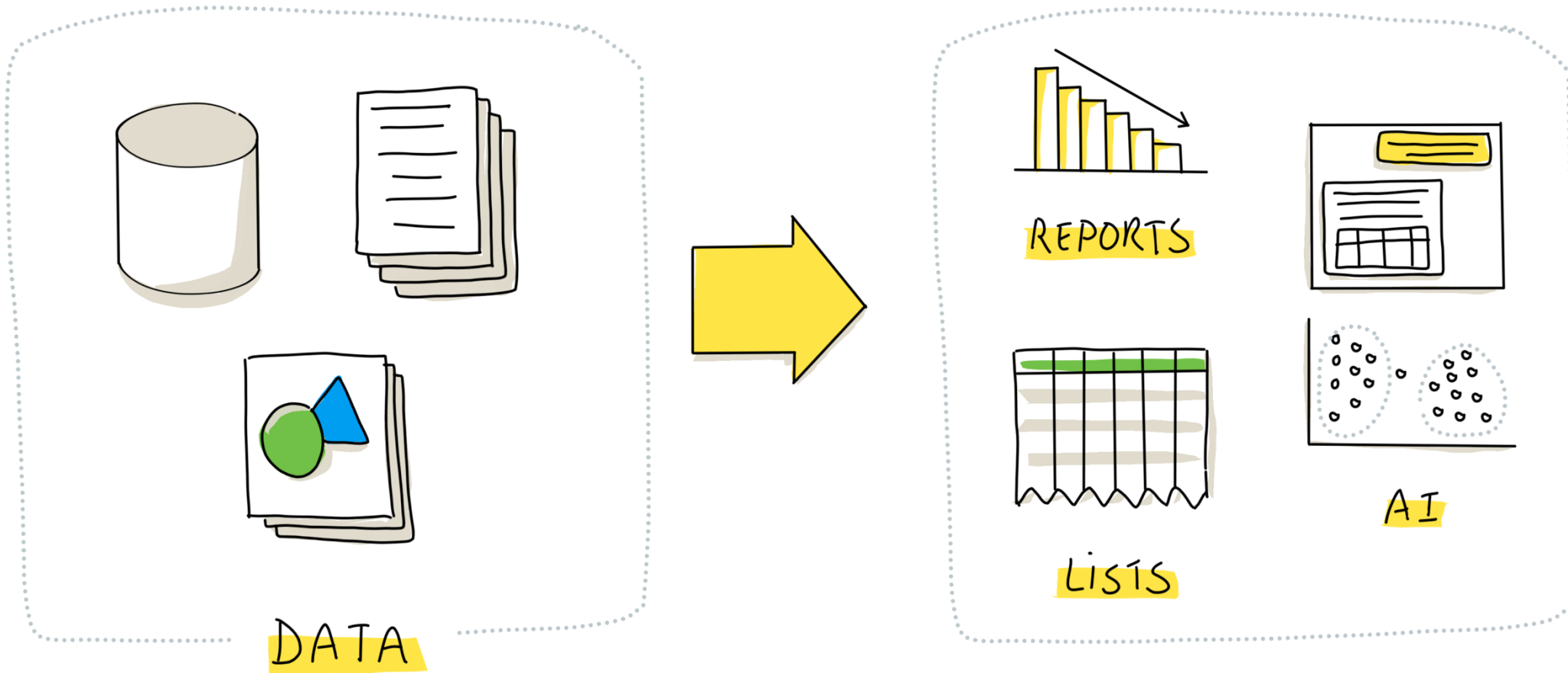


DEEL 9

DATA MANAGEMENT



Why Data Management in the context of AI?





Data Management

- Defining Data
- Data Producers
- Big Data Vs
- The Bigger Picture
- Data Management
- Data Technology



Data Management

- **Defining Data**
- Data Producers
- Big Data Vs
- The Bigger Picture
- Data Management
- Data Technology

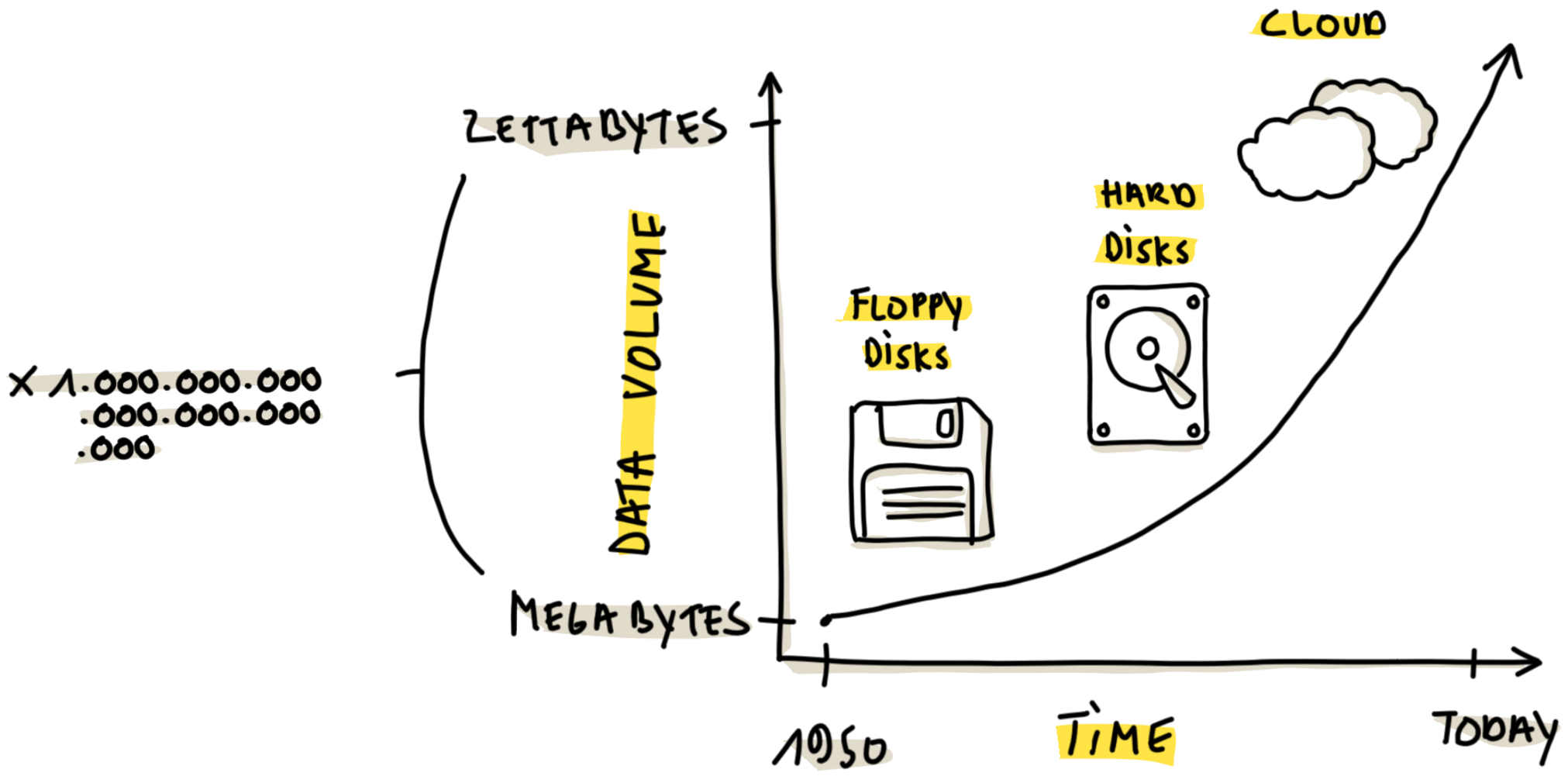


Defining Data

Data = Digital information that is **stored, processed, and manipulated** by computers and other devices.

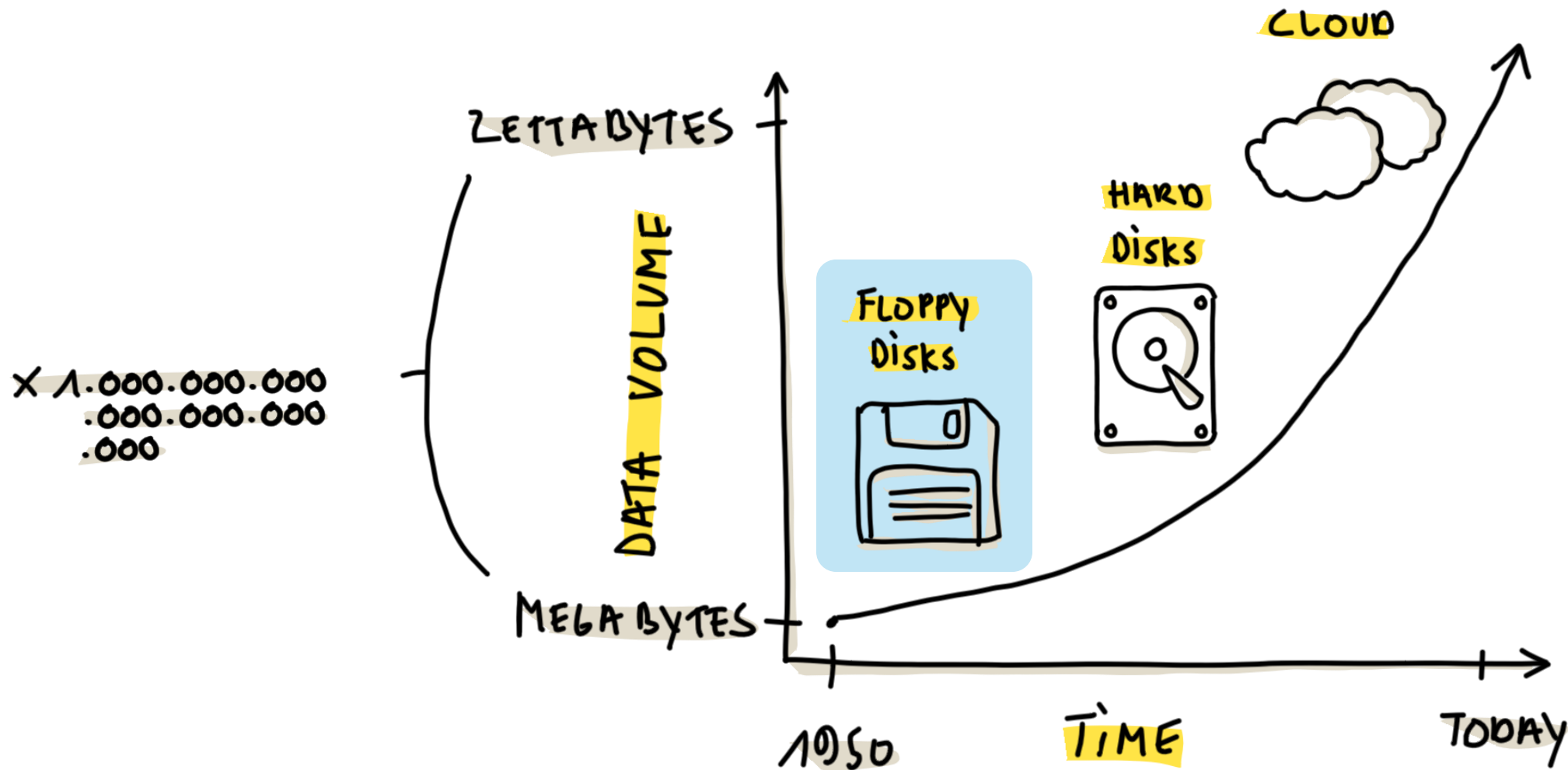
Examples of data: text, images, figures, video, spreadsheets, word documents, internet websites, sensor measurements, ...

The world of data changed a lot since the beginning (1950s) and today...





History of Data









720 KB



1.4-3 MB



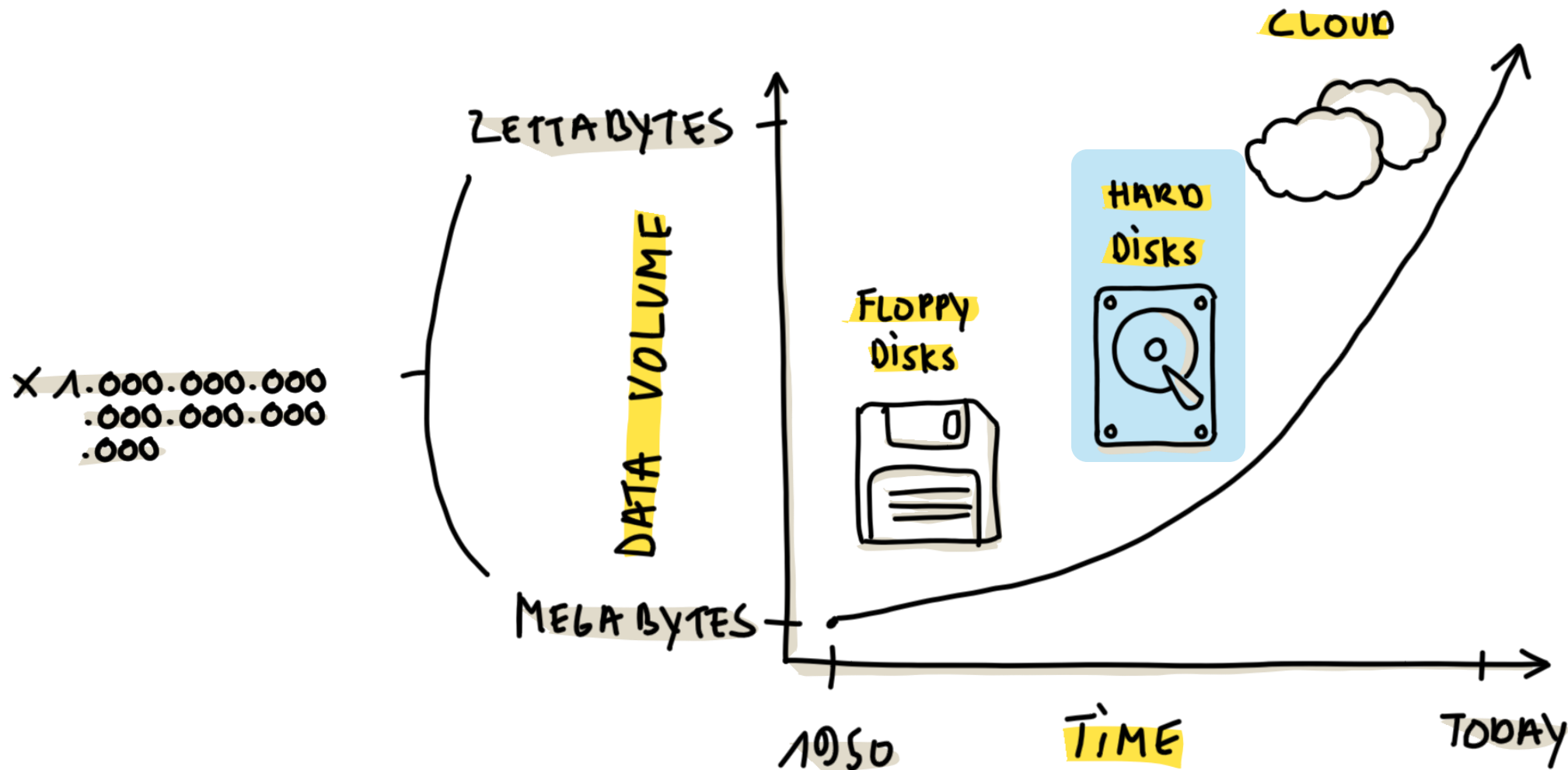


Floppy Disks

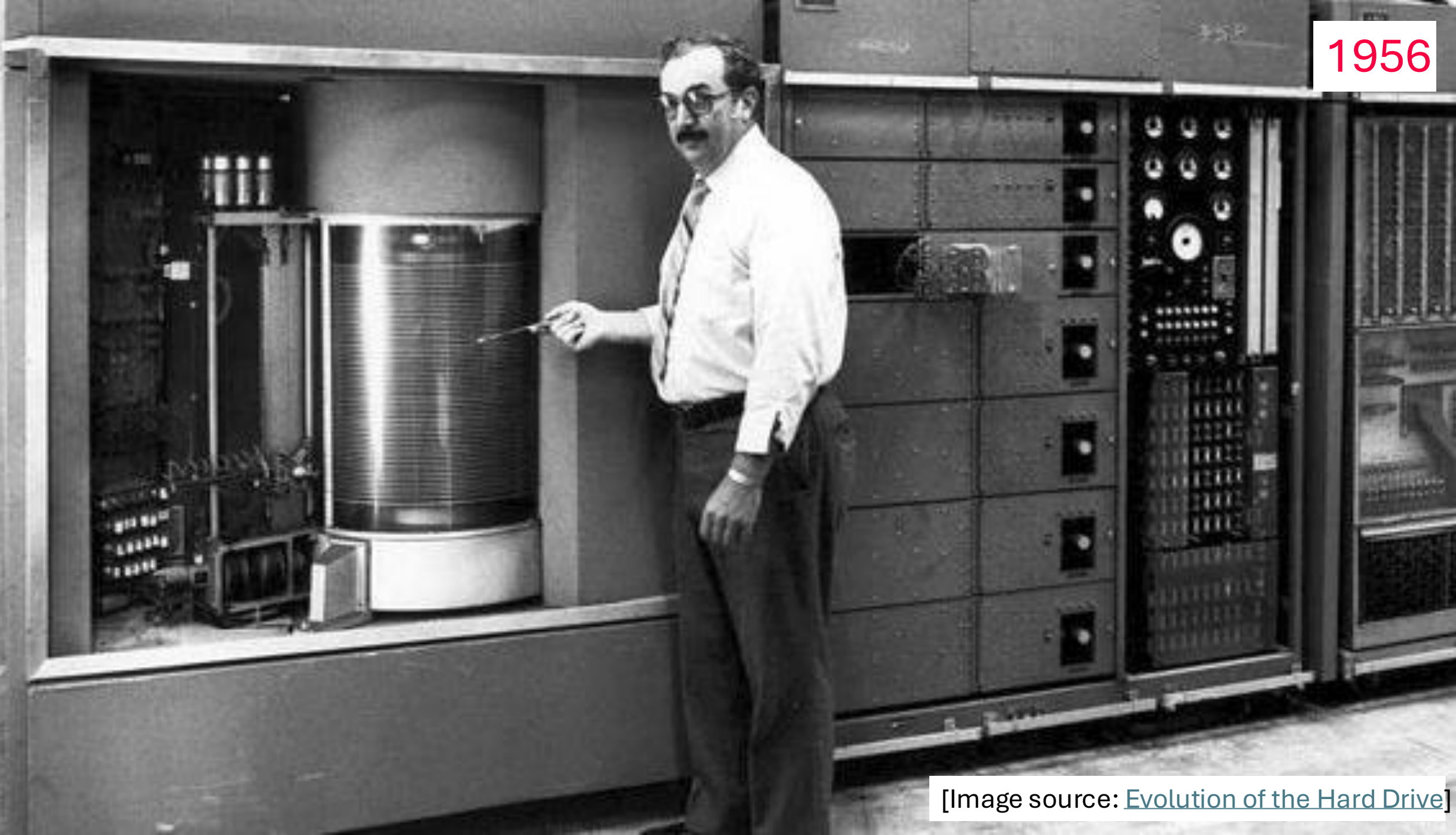
- Popular storage medium 1970s – end 1990s
- Size: kilobytes (KB) to a few megabytes (MB)
- In other words: storing a few word documents or a small image



History of Data



1956

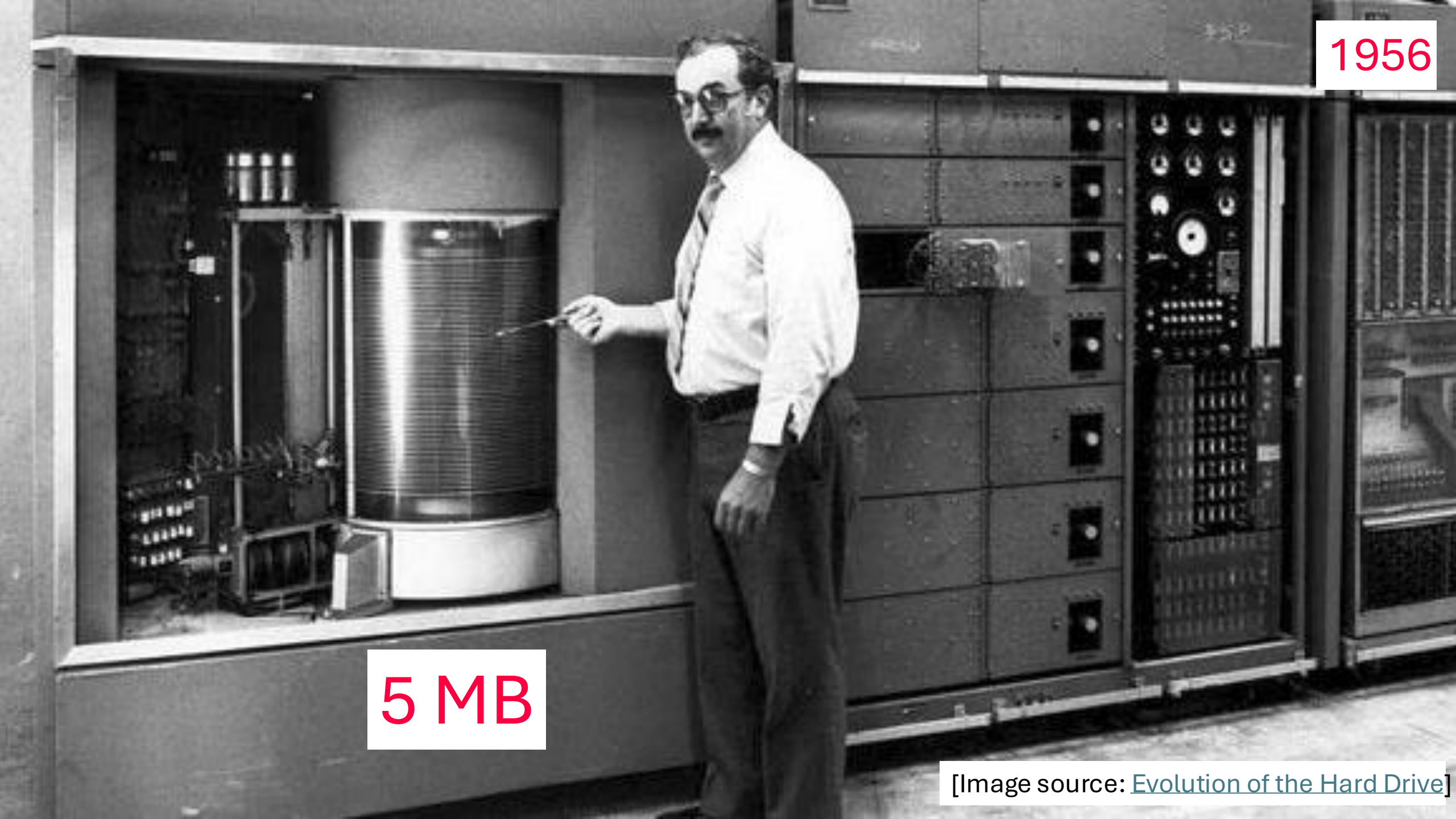


[Image source: [Evolution of the Hard Drive](#)]

1956

5 MB

[Image source: [Evolution of the Hard Drive](#)]



1970

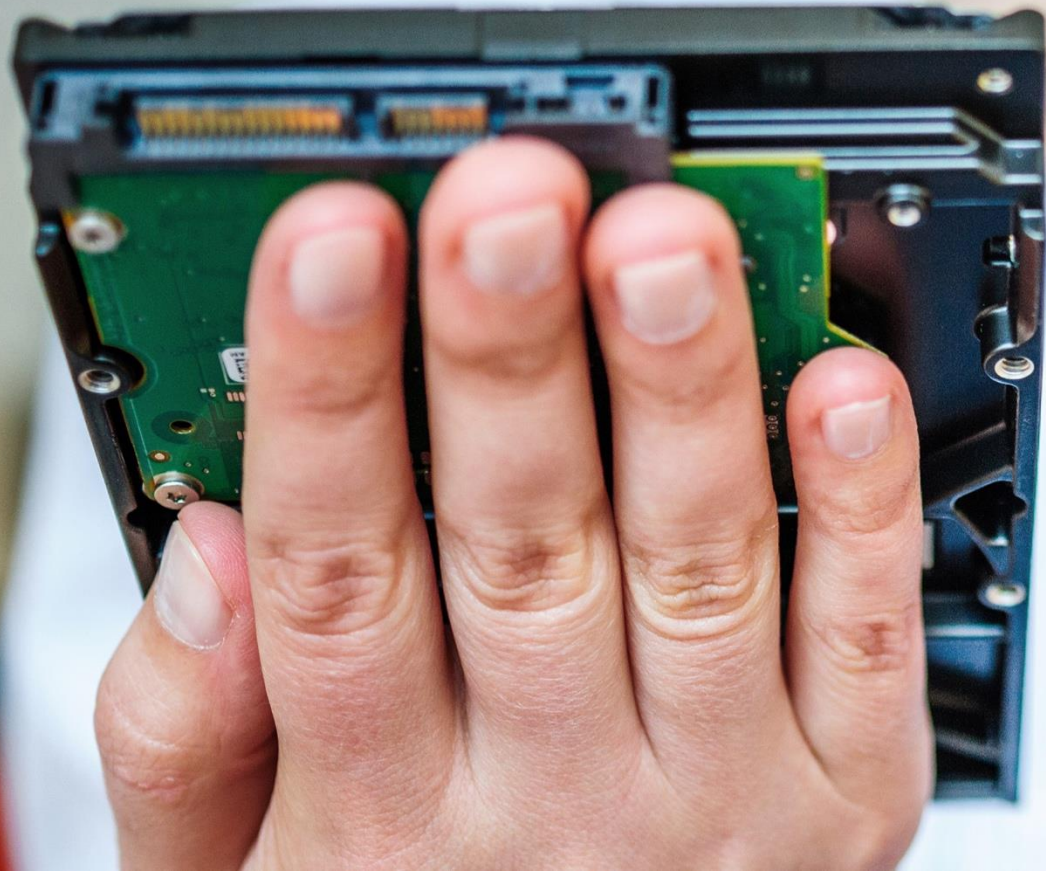


Removable Disk Packs

1980s - Today

HDD

SDD



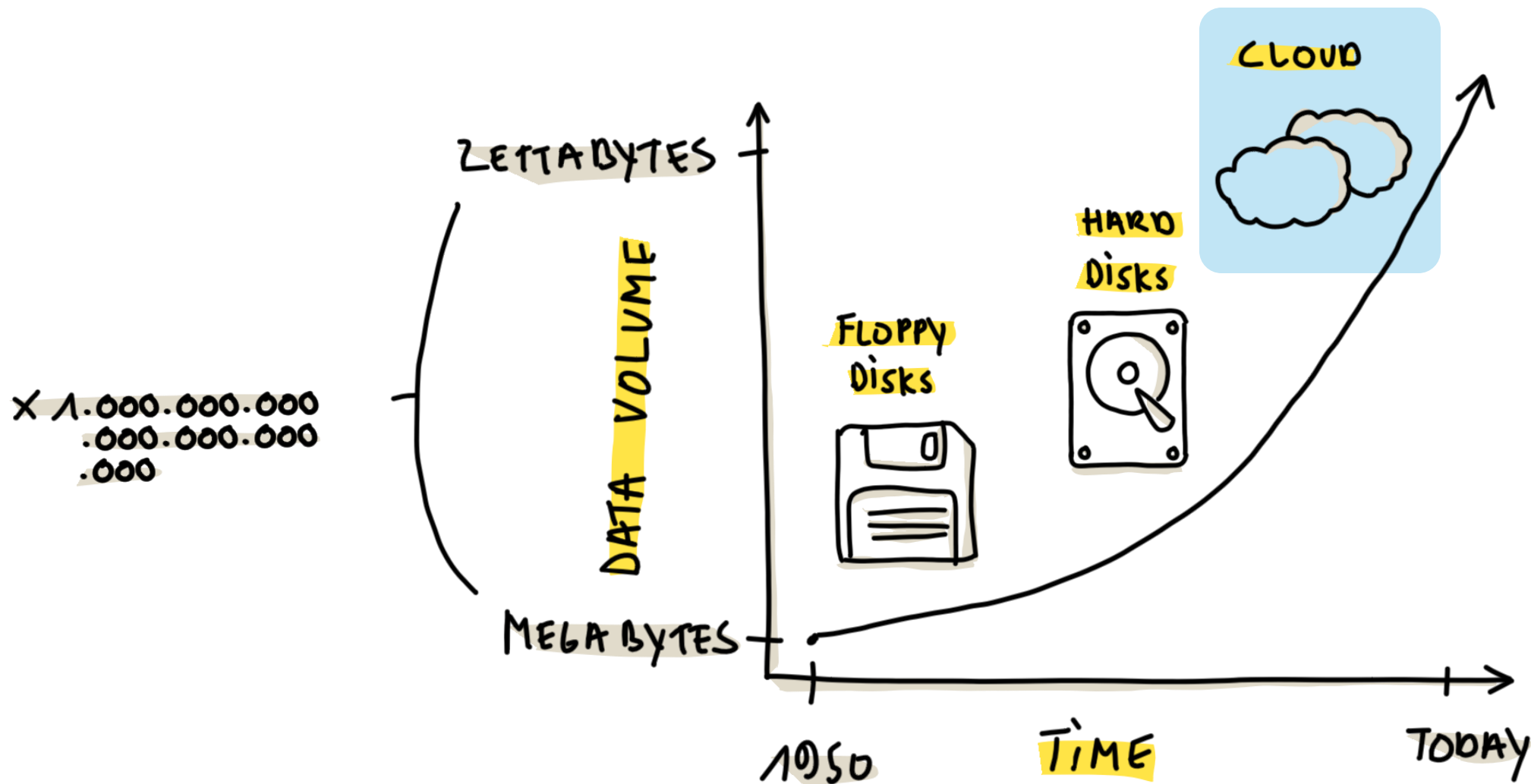


Hard Disk Drives (HDD)

- Introduced in 1950s
- Remains a key component of computer storage (datacenters, personal computers)
- Capacity and speed evolved over time
- Volume can range from Gigabytes (GB) to multiple Terrabytes (TB)
- 1 TB HDD = +/- 200.000 Photos
- Note: Smartphones, tablets, modern notebooks still have hard disks (Solid State Drives – SSDs)



History of Data





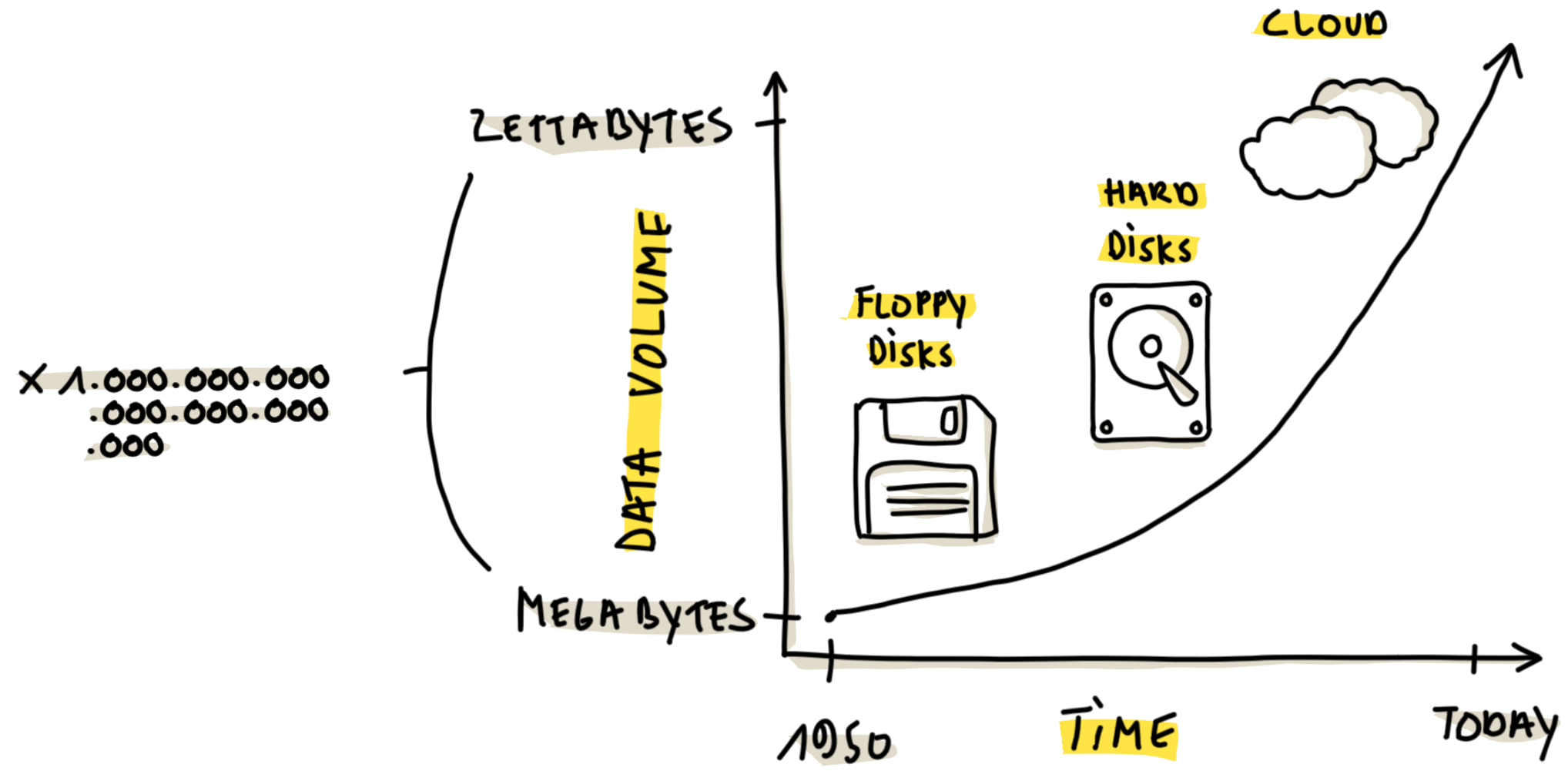
The Cloud



- Introduced in **the 2000s**
- Storing data on **remote servers** accessible via the Internet
- Cloud = **service based model**
- Example services Dropbox, Google Drive, iCloud, Amazon S3 or OneDrive
- **Upload, store and retrieve high volumes of data from anywhere**
- Cloud storage can **scale up or down** based on our needs
- **No more physical storage limits**

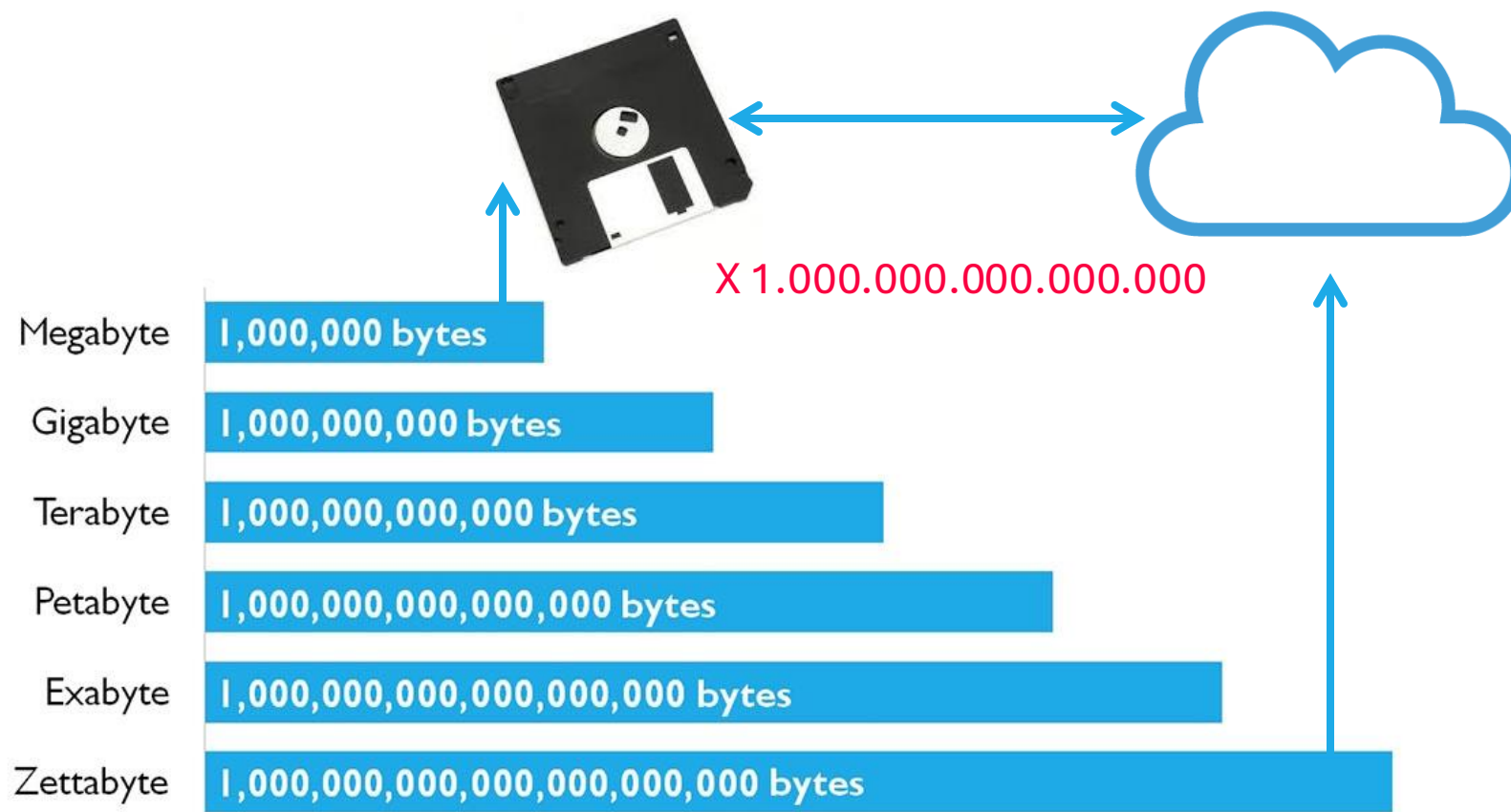


Some Observations...

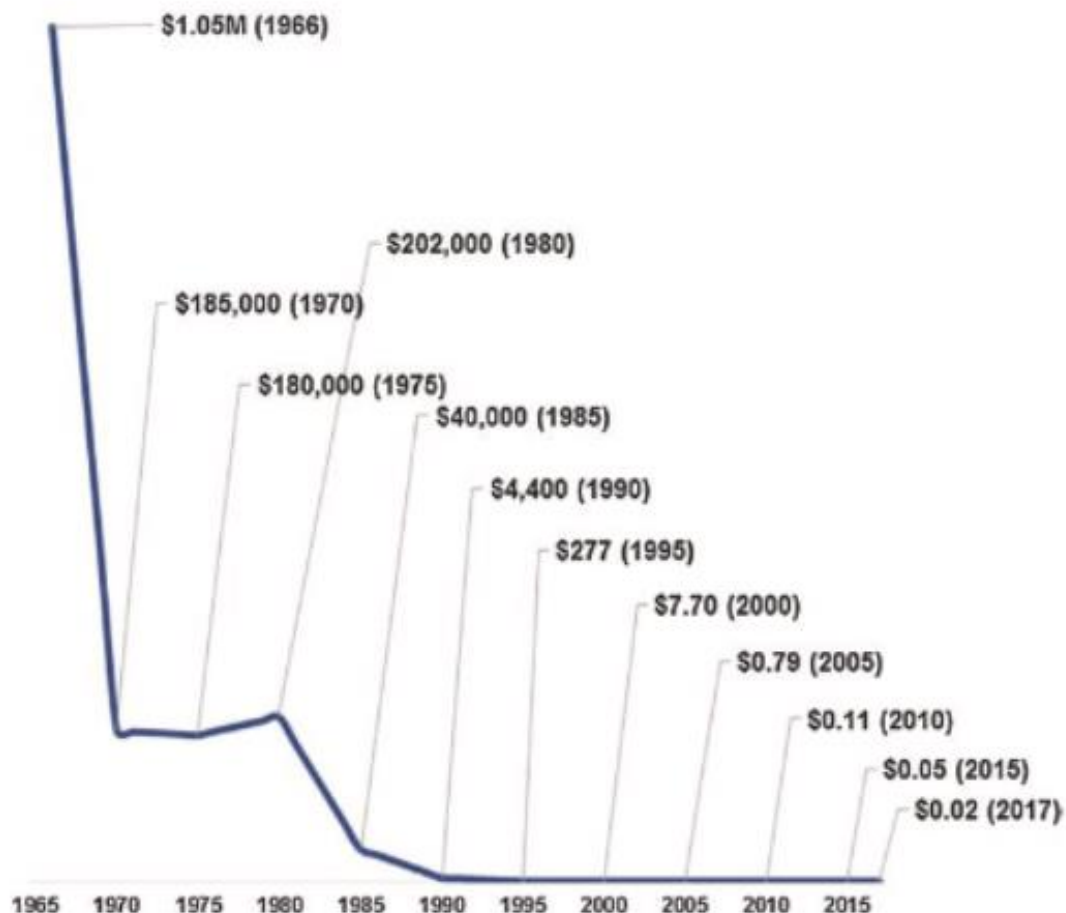




1) “No more physical storage limits”



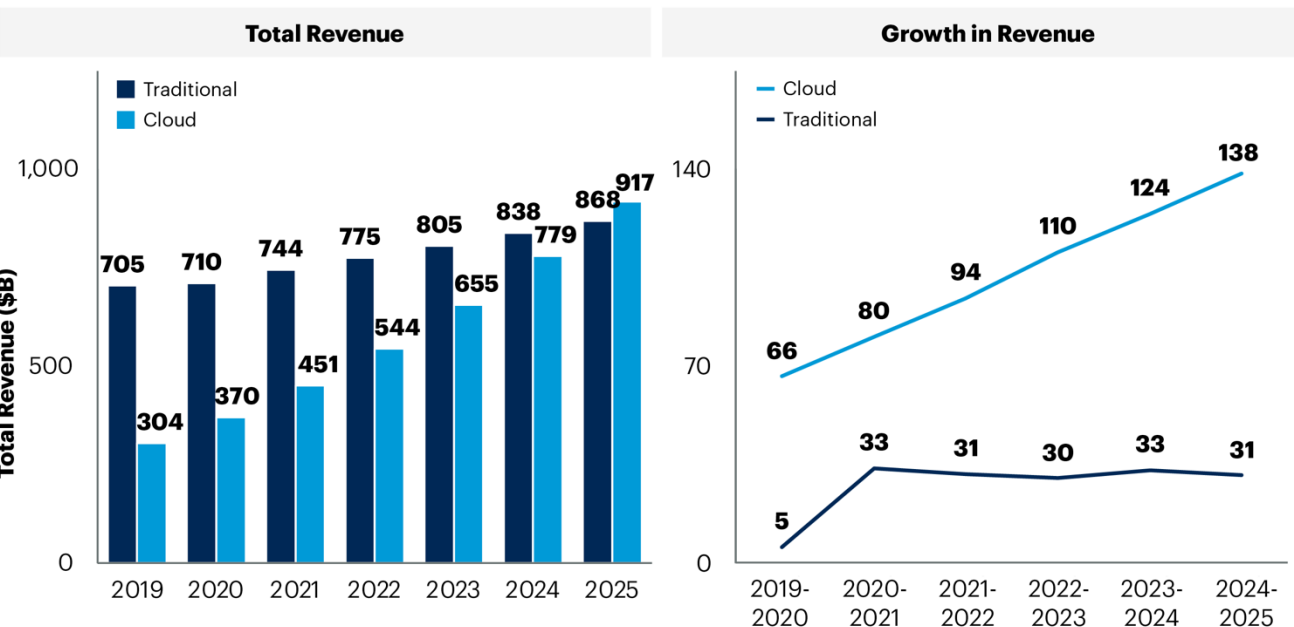
2) “Cost of storing data dropped”



- Graph: Cost per gigabyte of storage over time
- Once extraordinarily expensive, the cost of storing large quantities of data has dropped to fractions of a cent
- Source: [[The Future of Finance](#)]



3) “More and more enterprises move to the cloud”

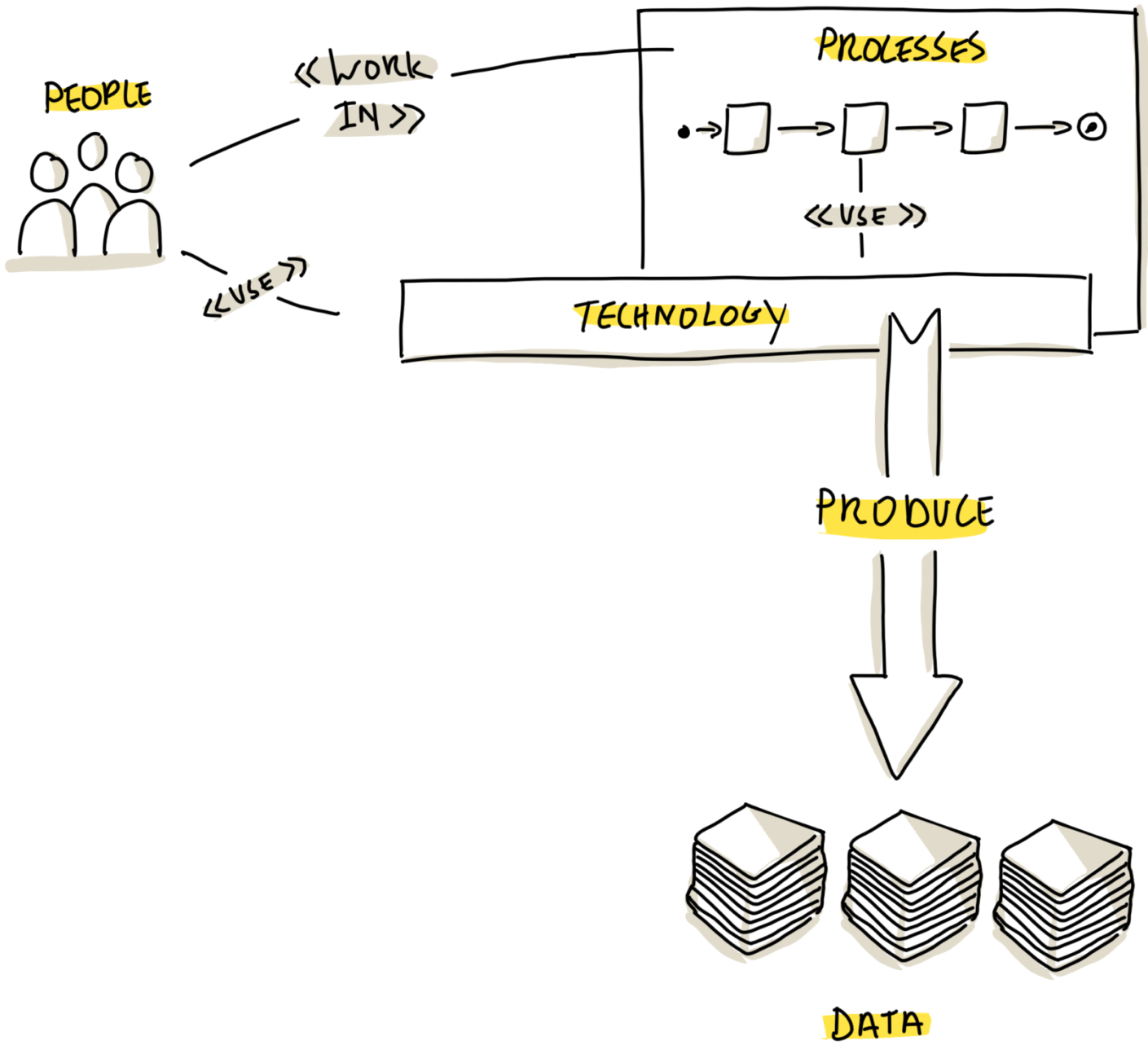


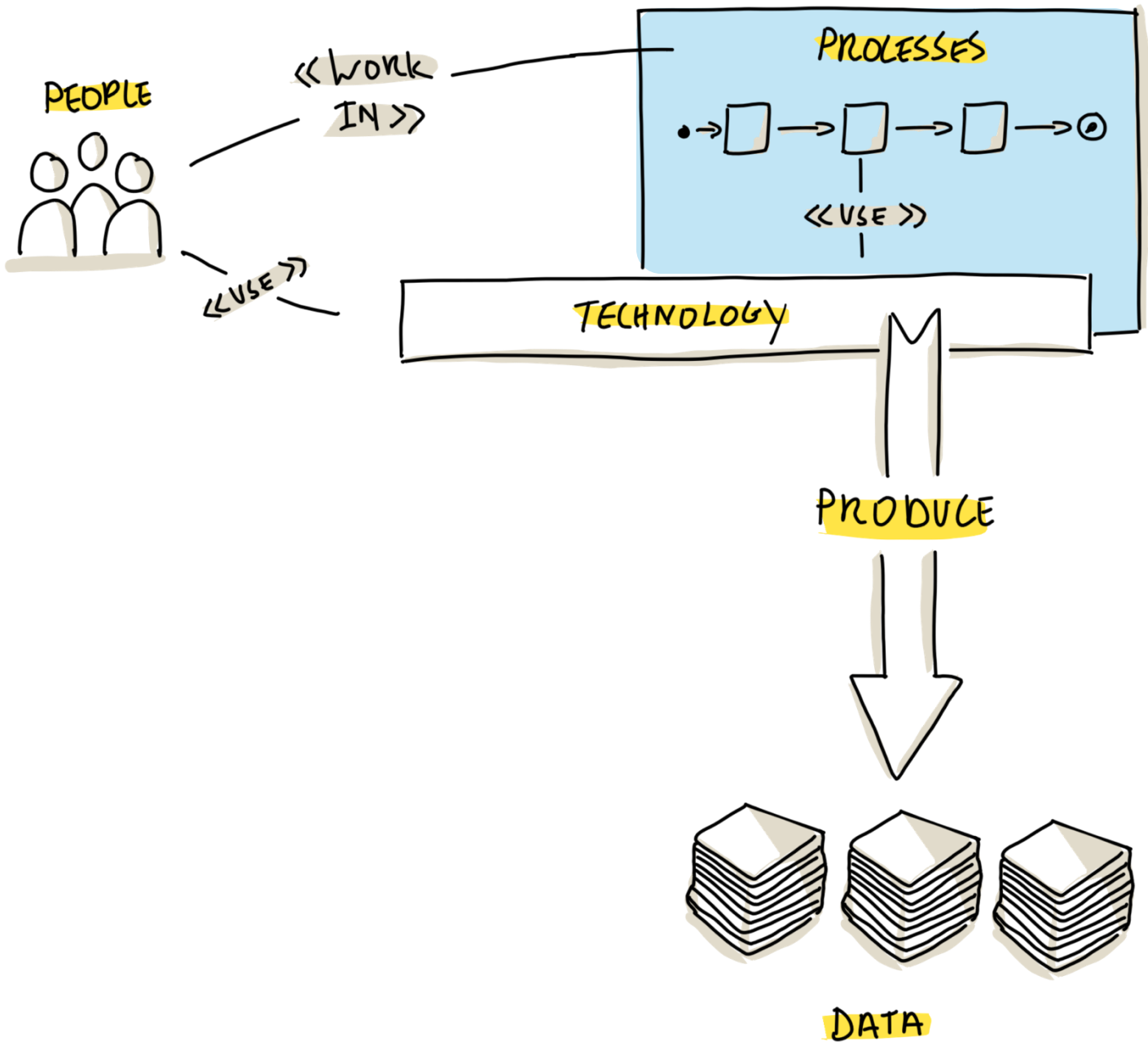
- Graph: Traditional vs Cloud computing
- **By 2025, 51% of IT spending will have shifted from traditional solutions to the cloud, compared to 41% in 2022.**
- Source: [[Gartner Press Release](#)]
- Cloud is an opportunity to work with higher volumes of data



Data Management

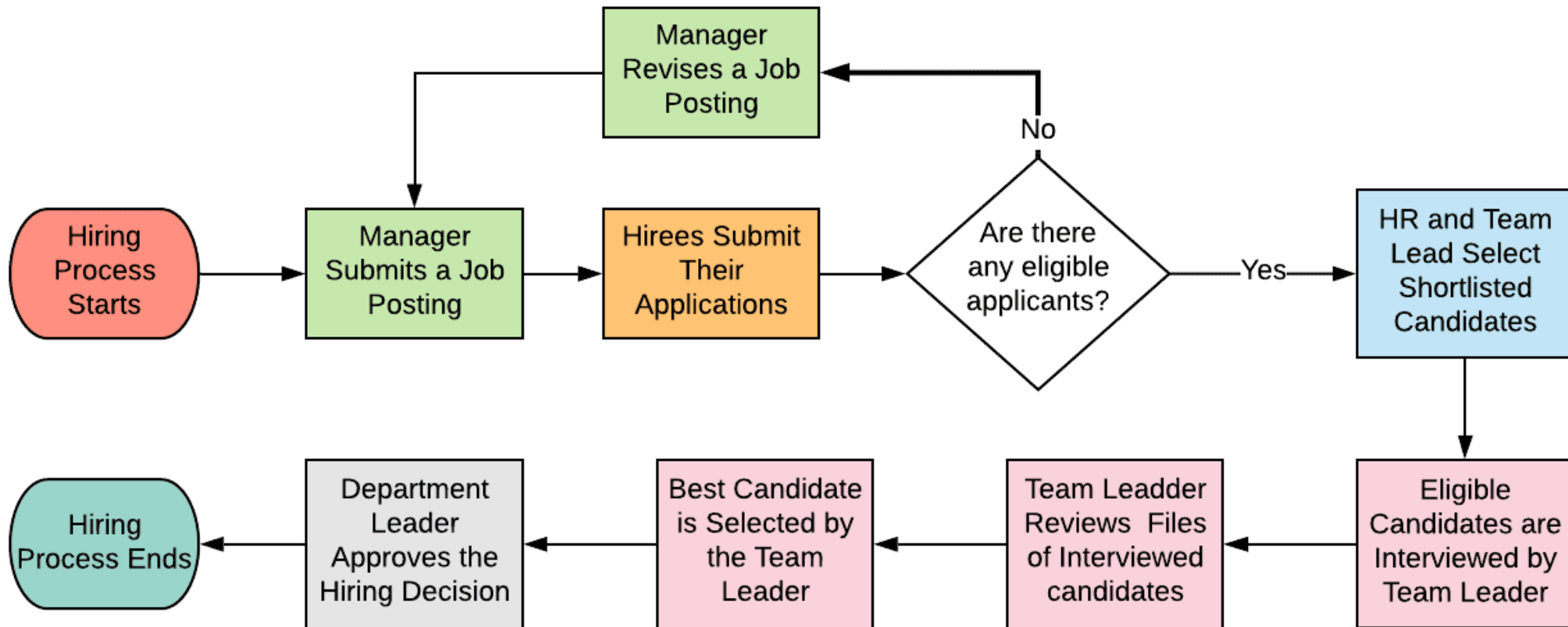
- Defining Data
- **Data Producers**
- Big Data Vs
- The Bigger Picture
- Data Management
- Data Technology







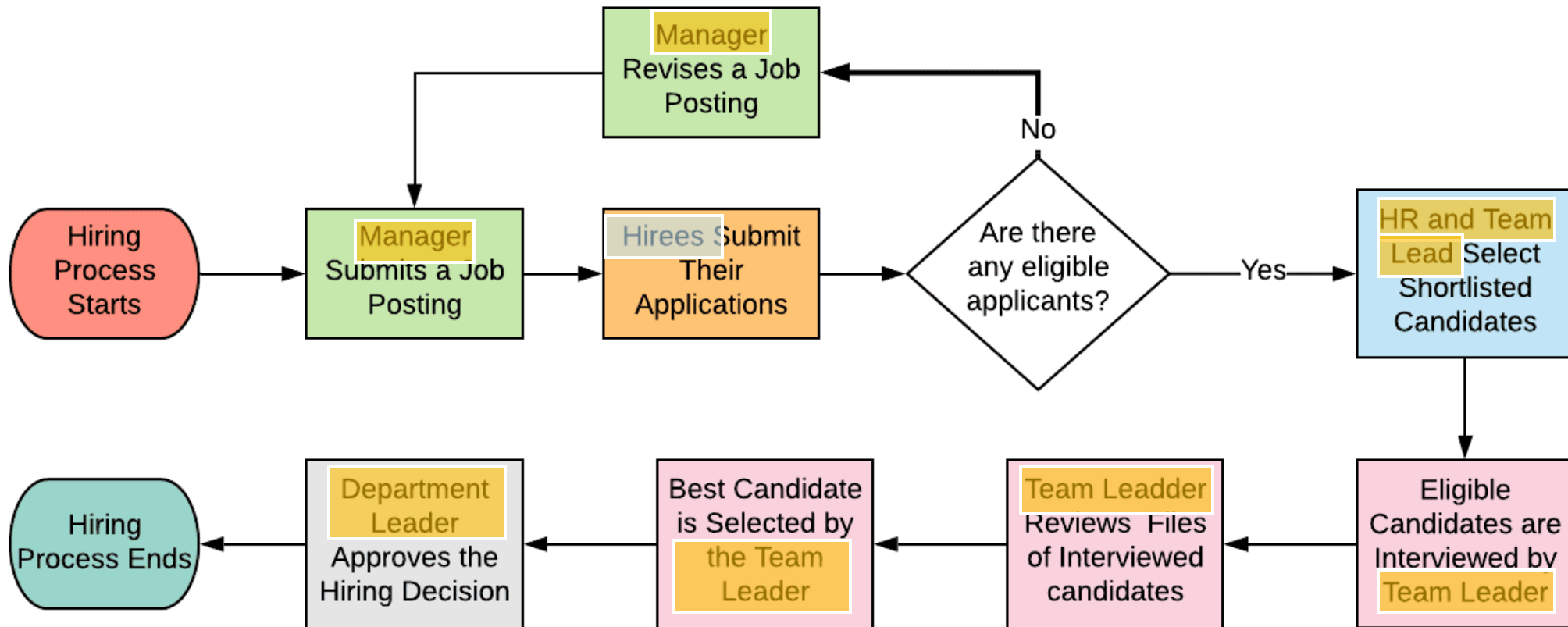
Process Example: Hiring a new Employee





Process Example: Hiring a new Employee

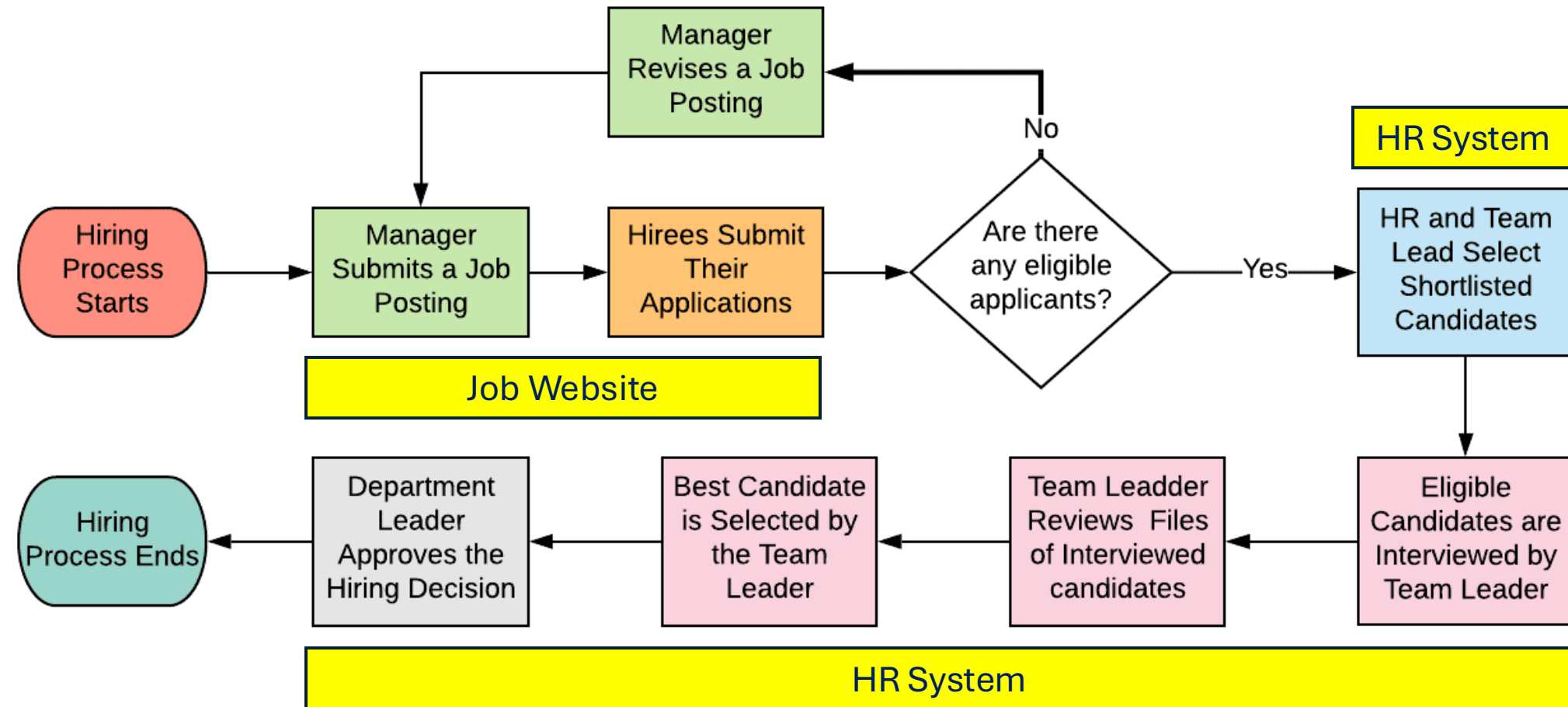
PEOPLE WORK IN THE PROCESS

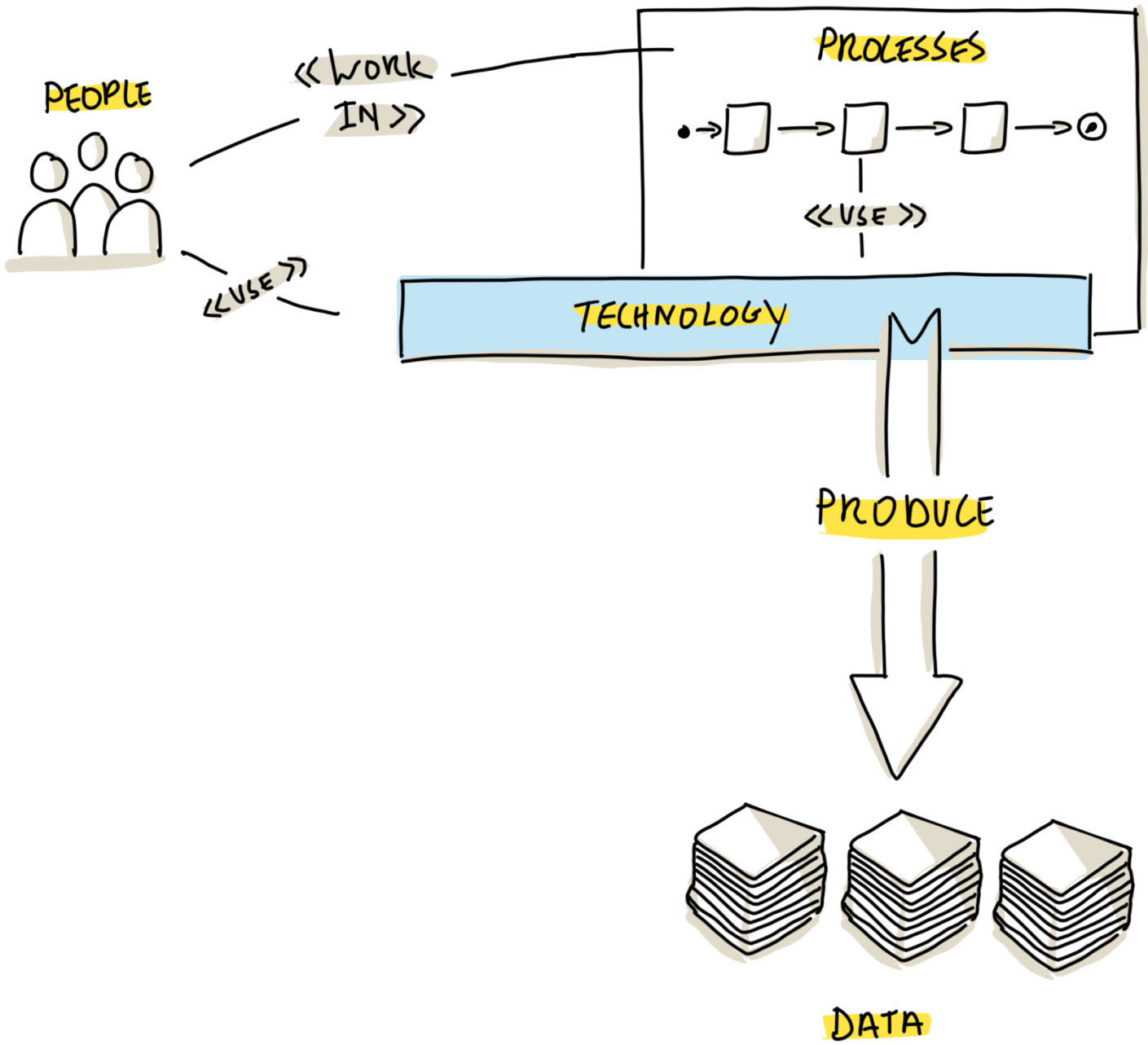


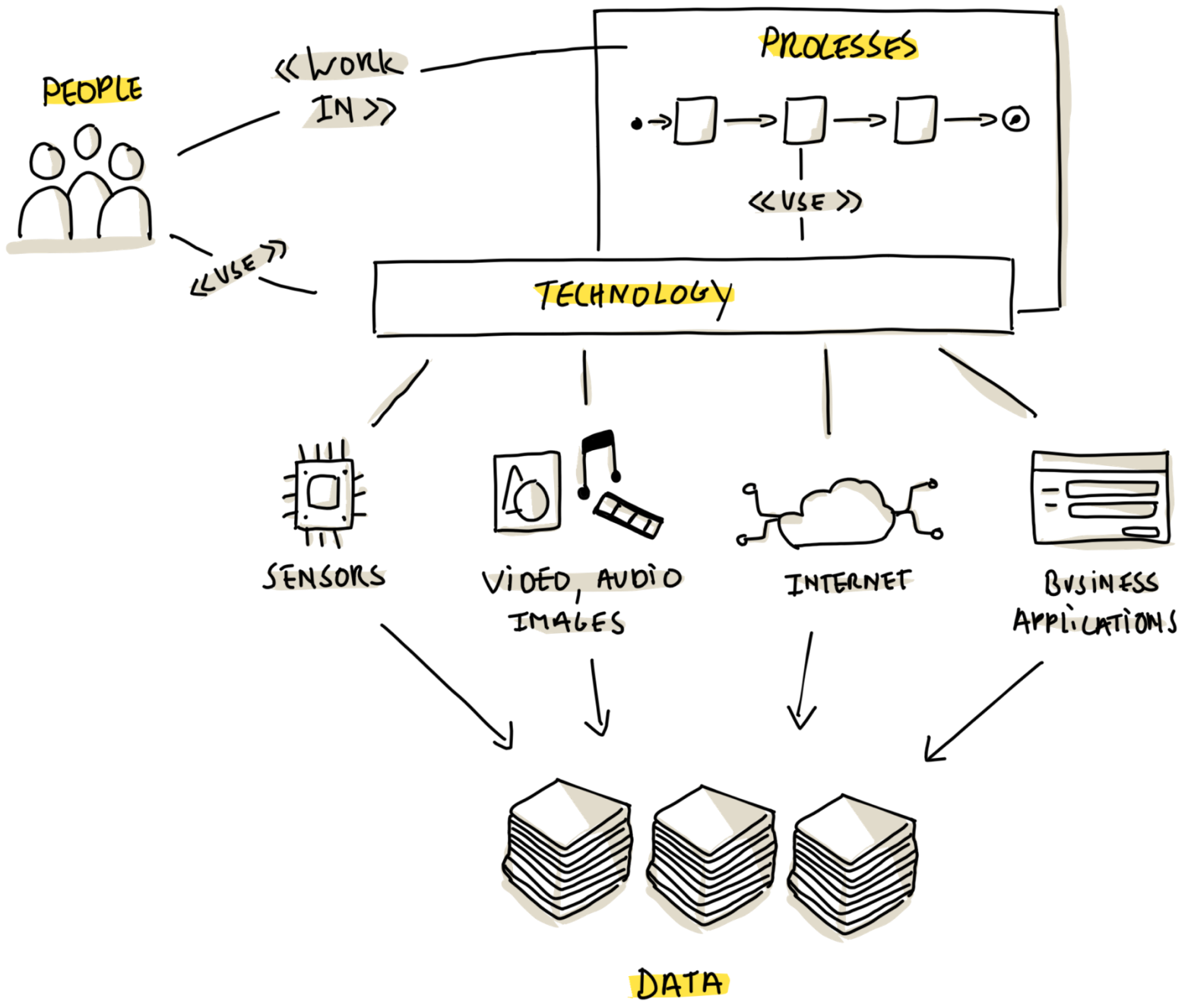


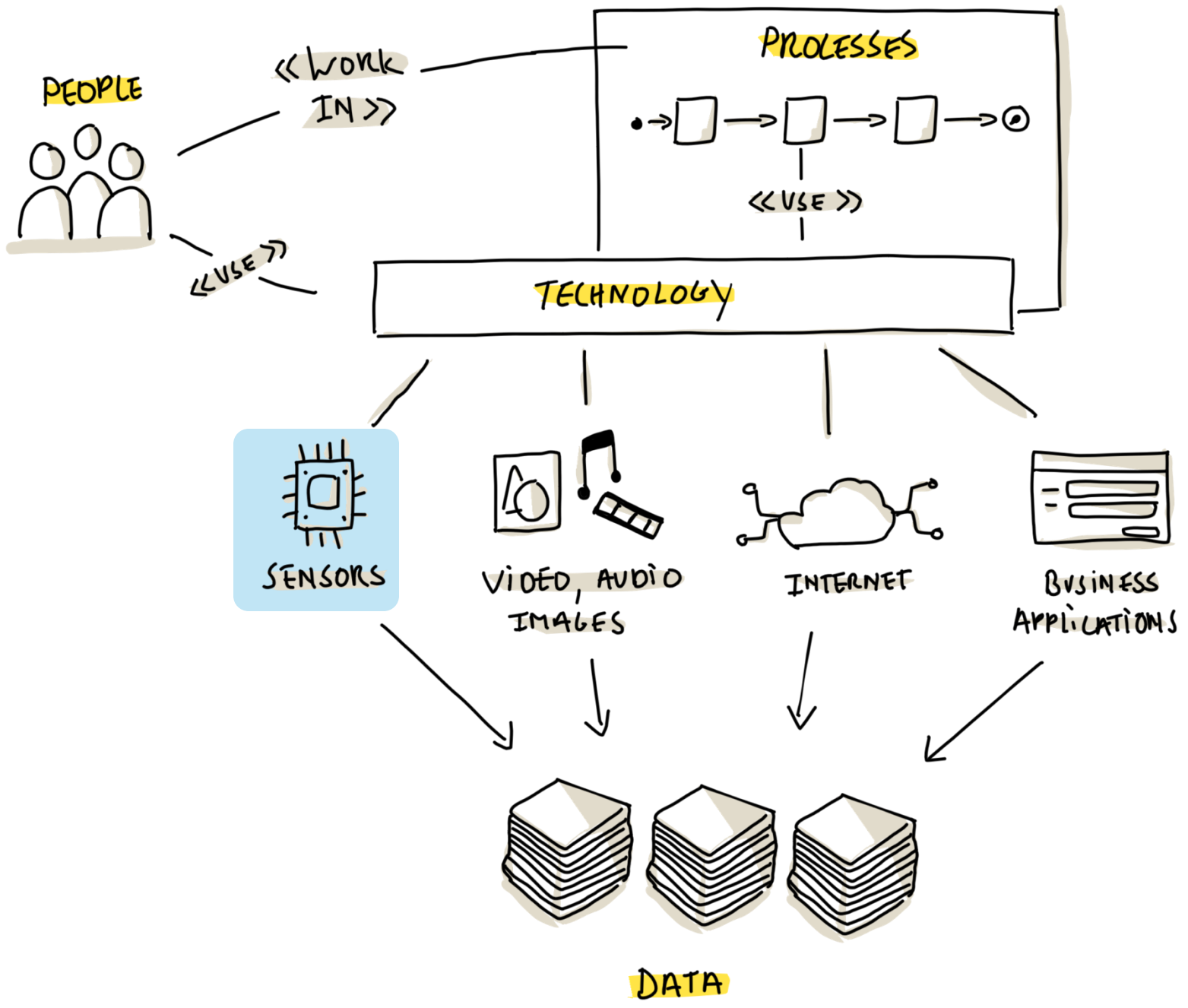
Process Example: Hiring a new Employee

TECHNOLOGY IS BEING USED IN THE PROCESS



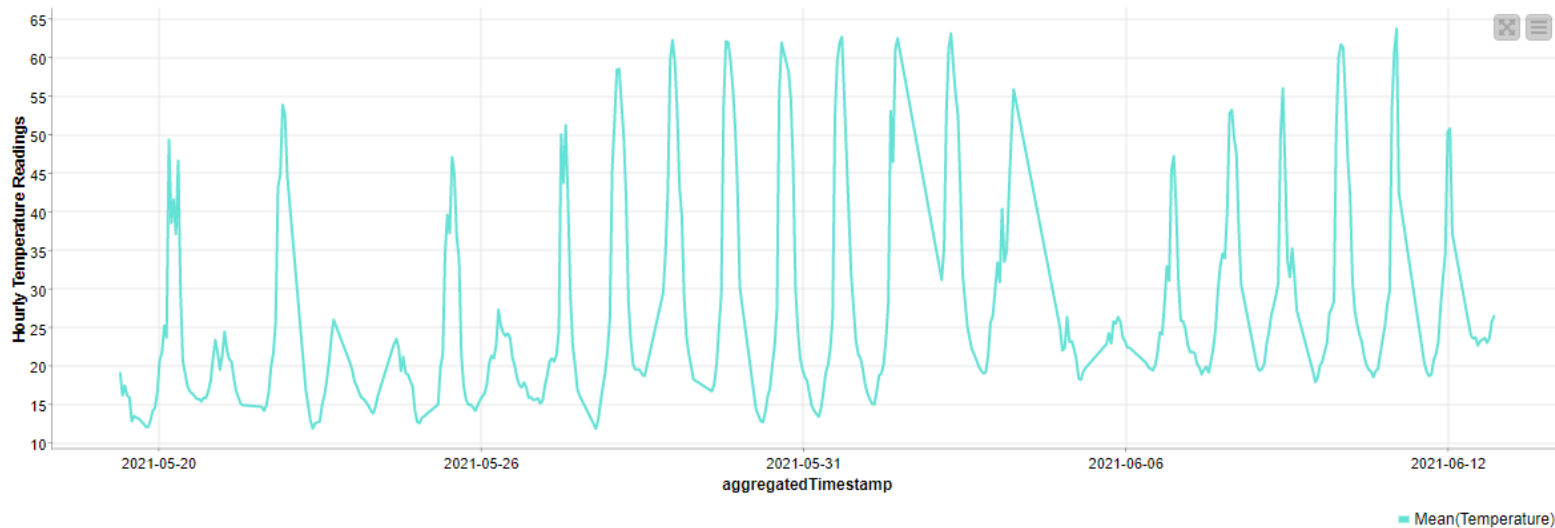








Temperature Sensors



Every 10 milliseconds, a temperature sensor will measure temperature.
This gives **8.640.000** temperature values per day.

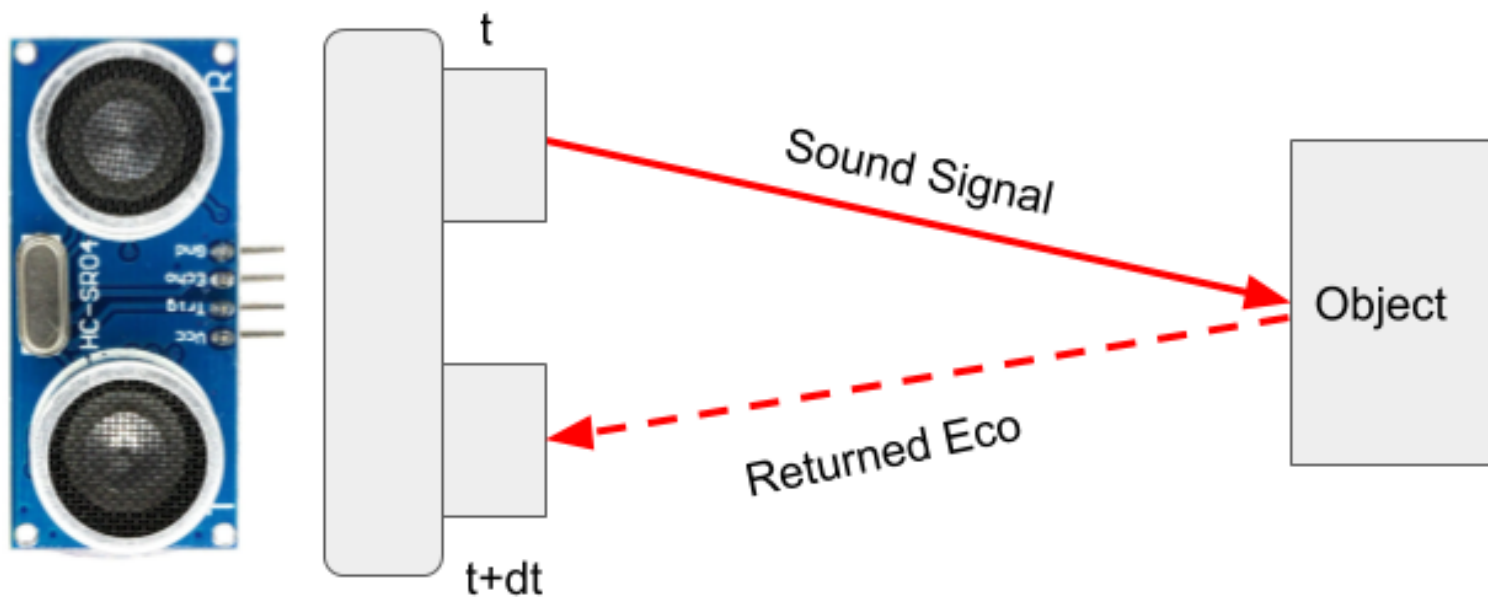


Temperature Sensors: Usage

- Controlling the temperature throughout process steps
- Example applications:
 - Food production
 - Chemical processes
 - Agriculture
 - Heating control
 - Transport Refrigeration
 - ...



Proximity Sensors

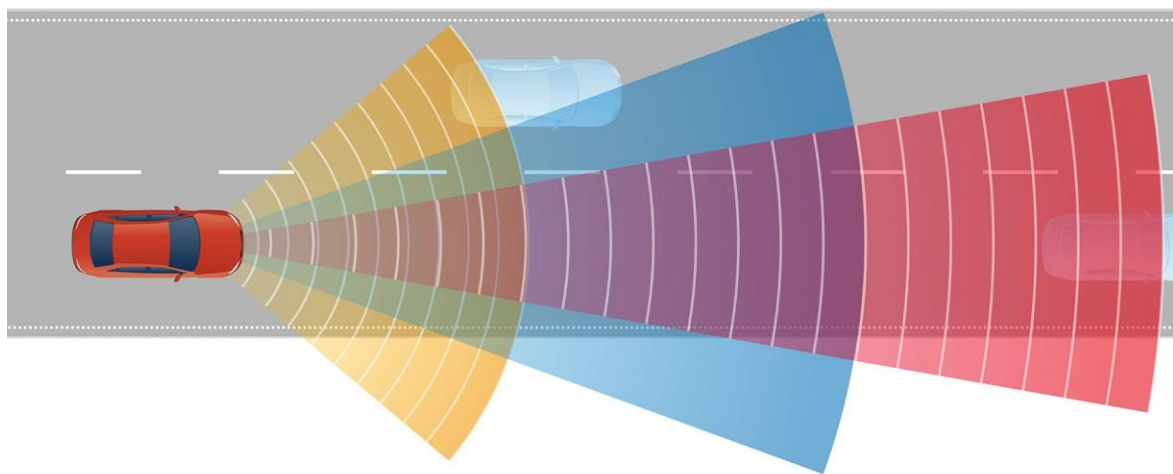


Proximity sensors detect the **presence or absence of a nearby object** or material.



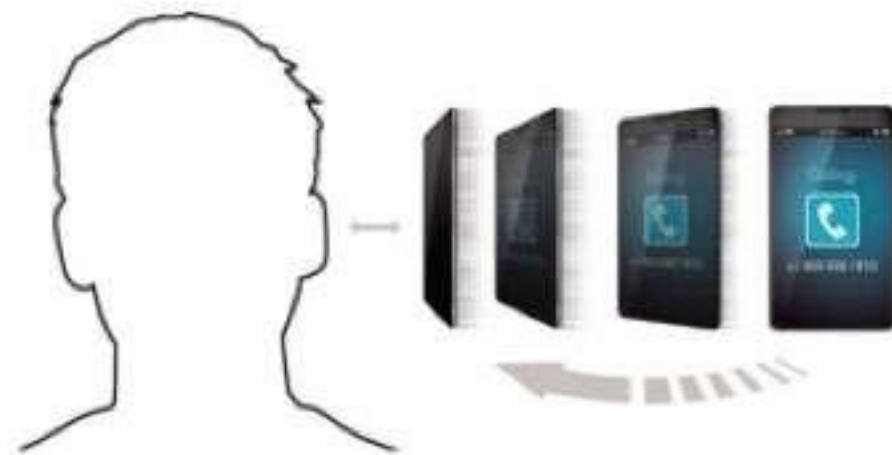
Proximity Sensors: Cars & Smartphones

CARS



Proximity sensors detect nearby objects to facilitate driving assistance.

SMARTPHONES

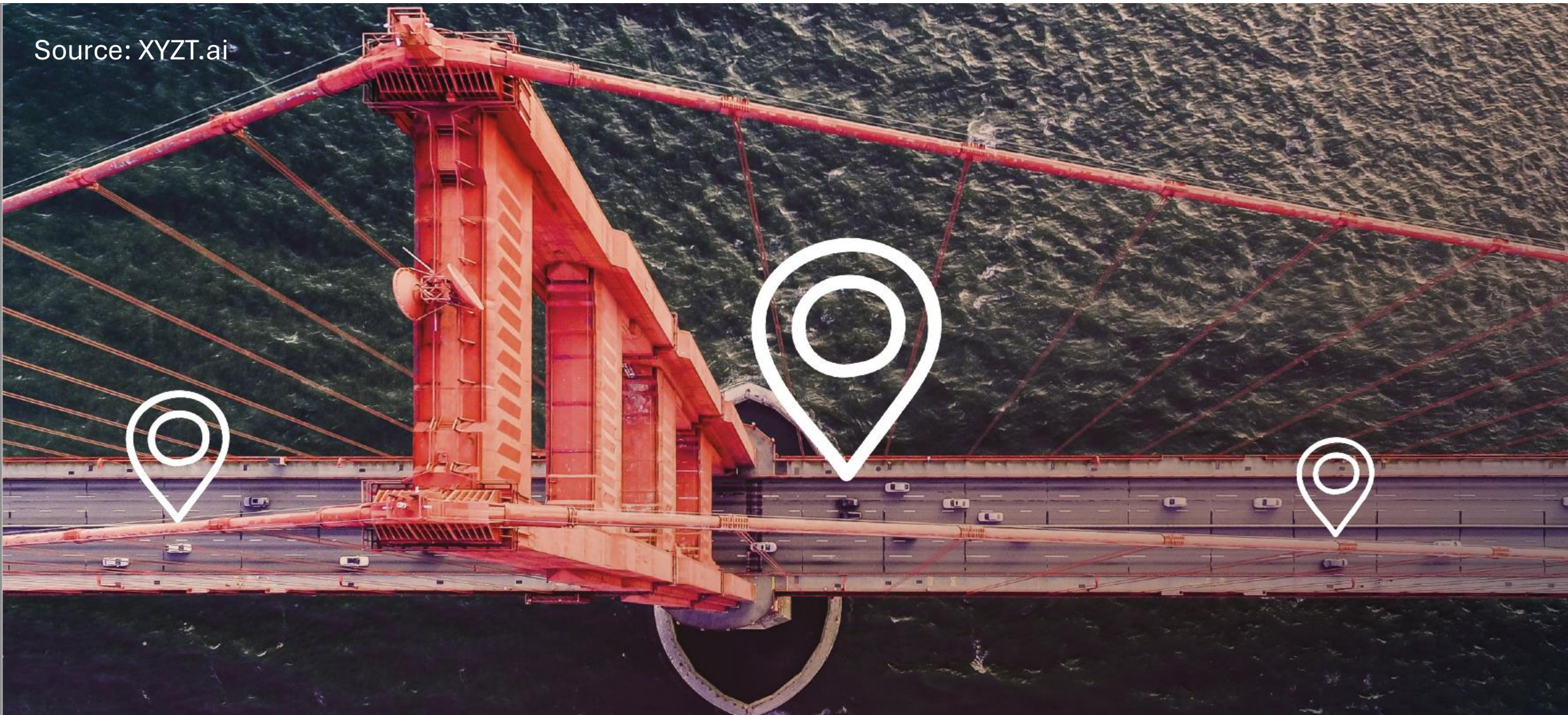


When a proximity sensor detects a face/ear, the smartphone screen turns off (energy saving).



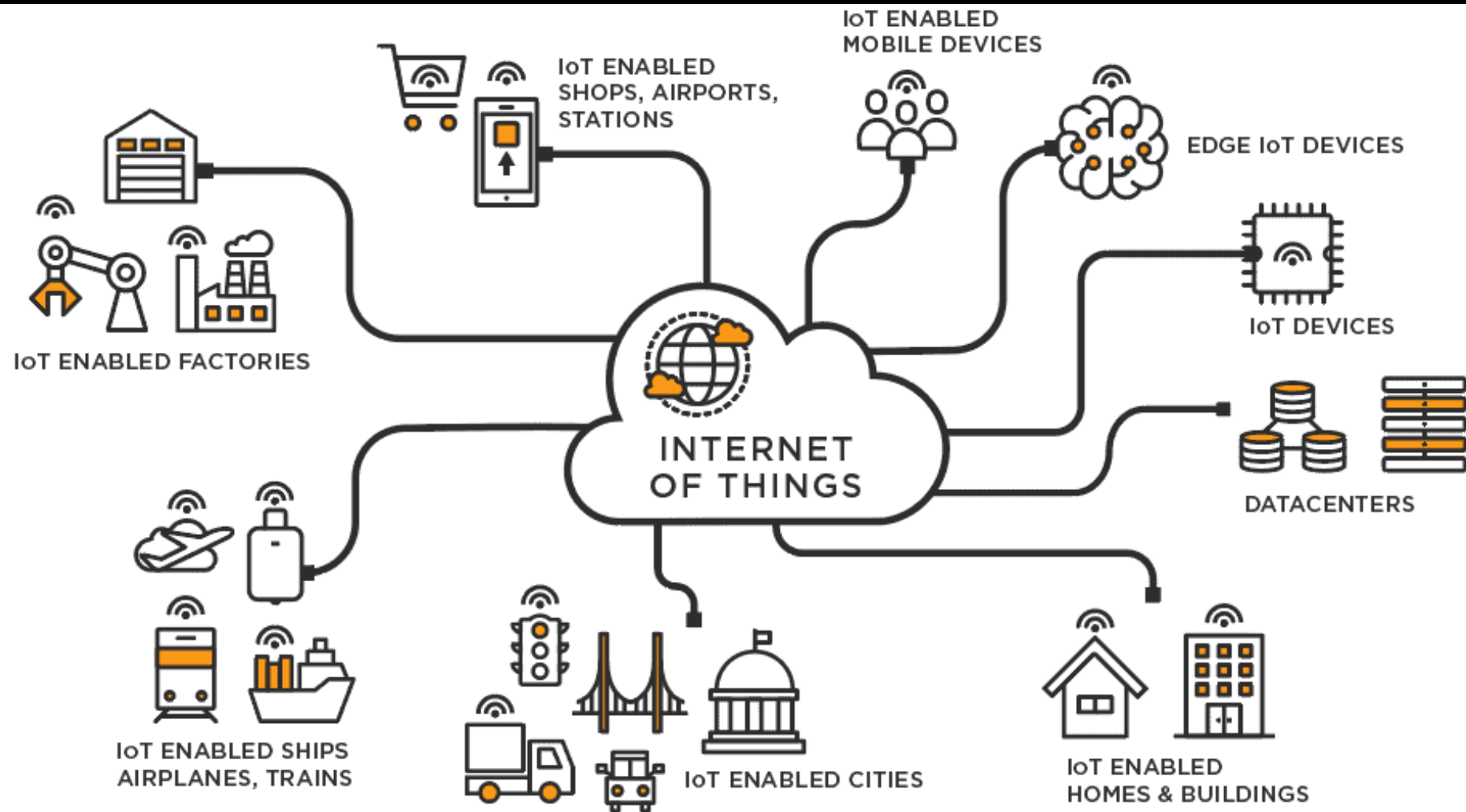
GPS Location Sensors

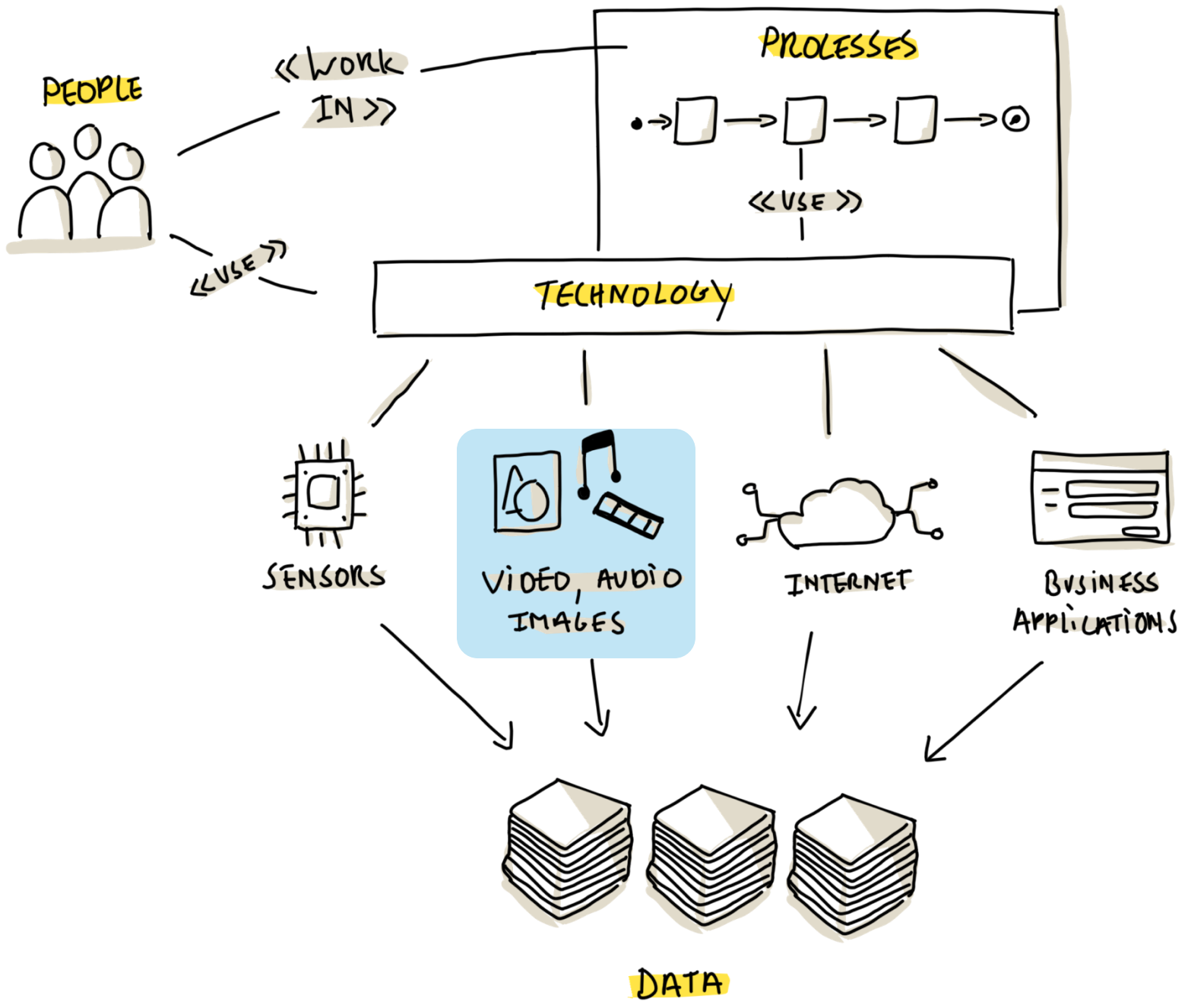
Source: XYZT.ai





Big Data Producers: Internet Of Things (IoT)





Video

- Data generated by cameras
- Various types:
 - **Regular cameras**, just recording video
 - **Smart cameras**, recording video + object/precense recognition
 - **Sattelite** cameras
 - **Drone** cameras





Smart Camera's with Object Detection





Smart Camera's with Object Detection





Volume

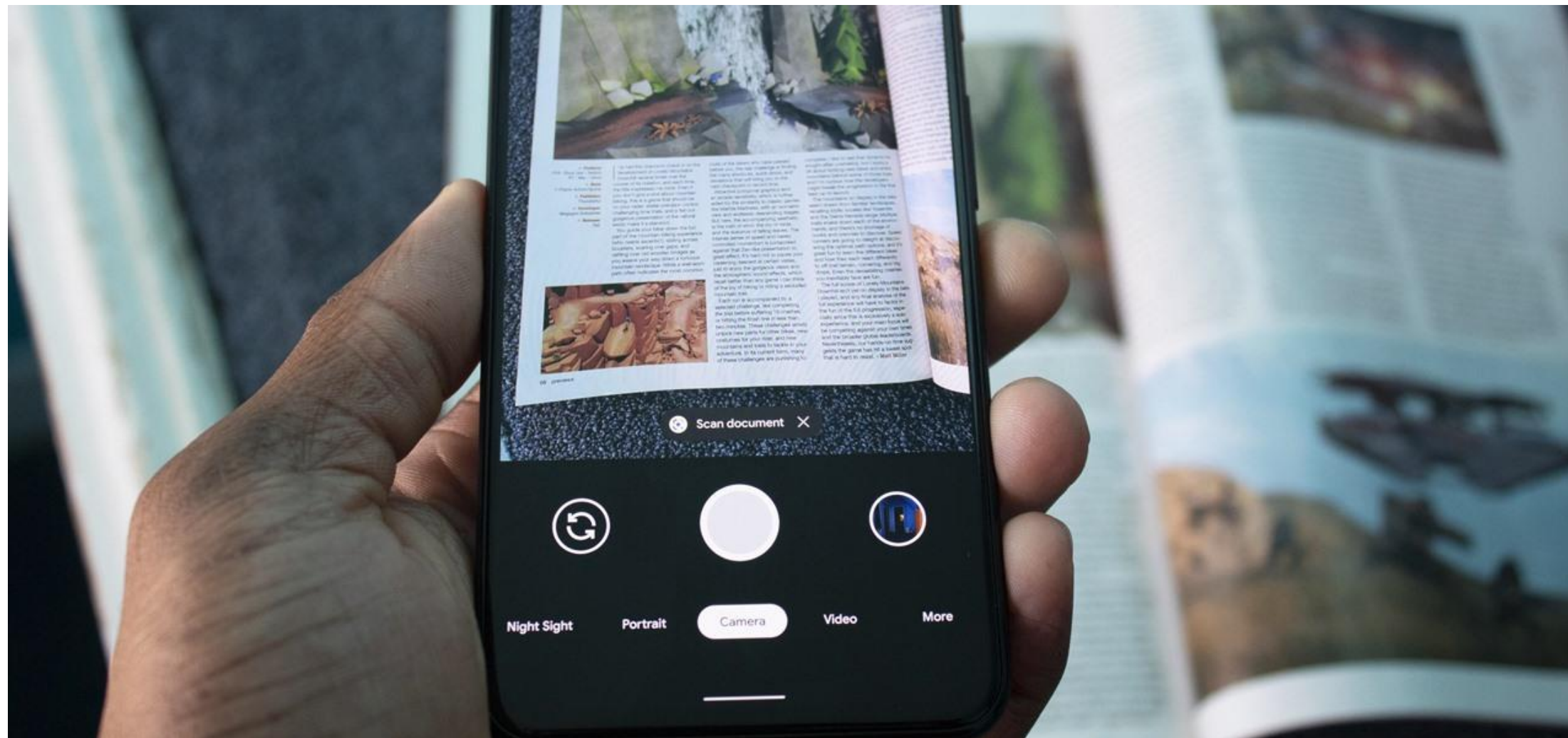
← **LOW** VIDEO QUALITY **HIGH** →

	144p	240p	360p	480p	720p	1080p
Per minute	1.3MB	3.3MB	5MB	8.3MB	25MB	50MB
Per hour	80MB	200MB	300MB	500MB	1.5GB	3GB

1 minute/hour of video data is **way larger** than 1 minute/hour of sensor data.

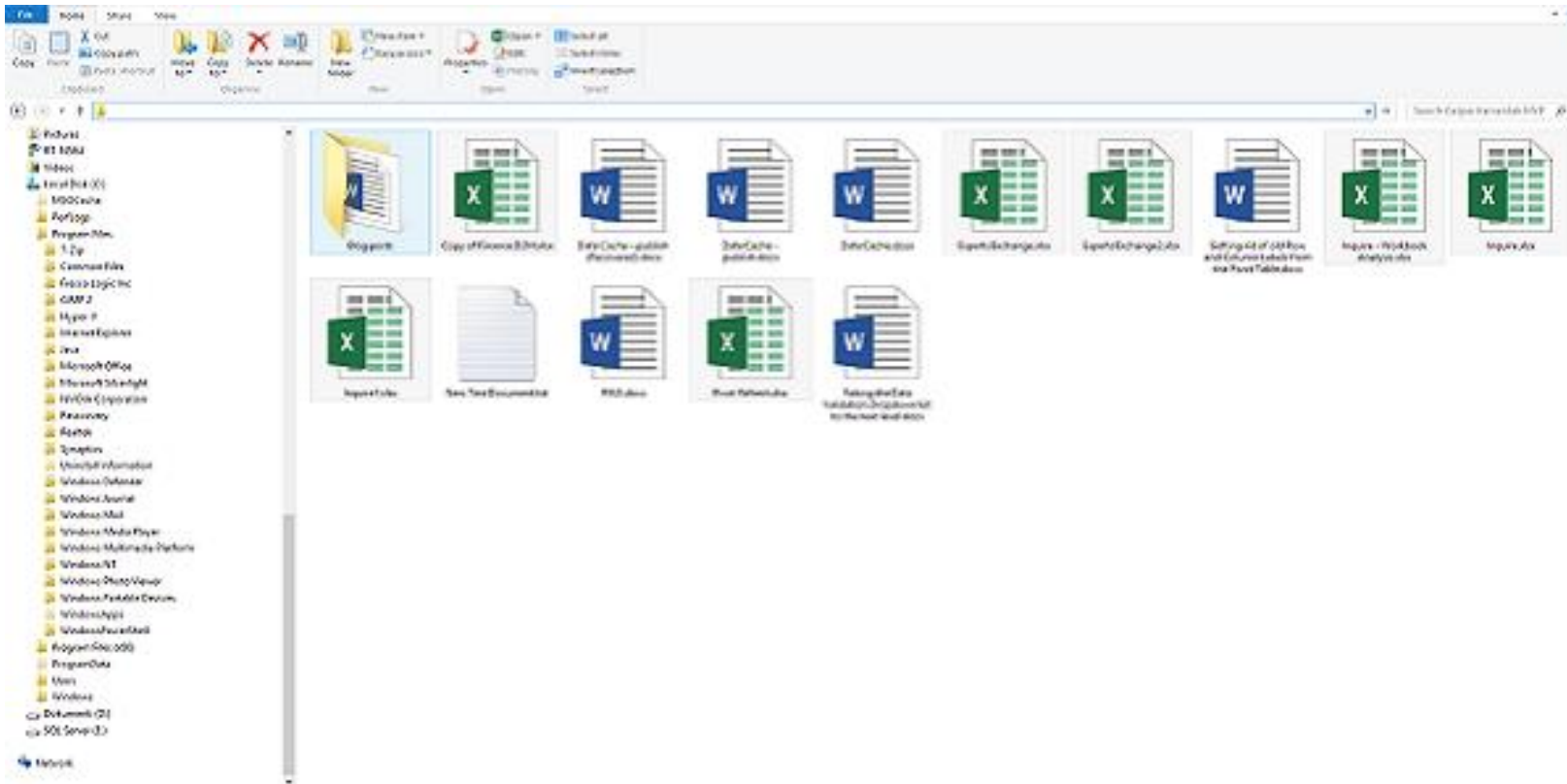


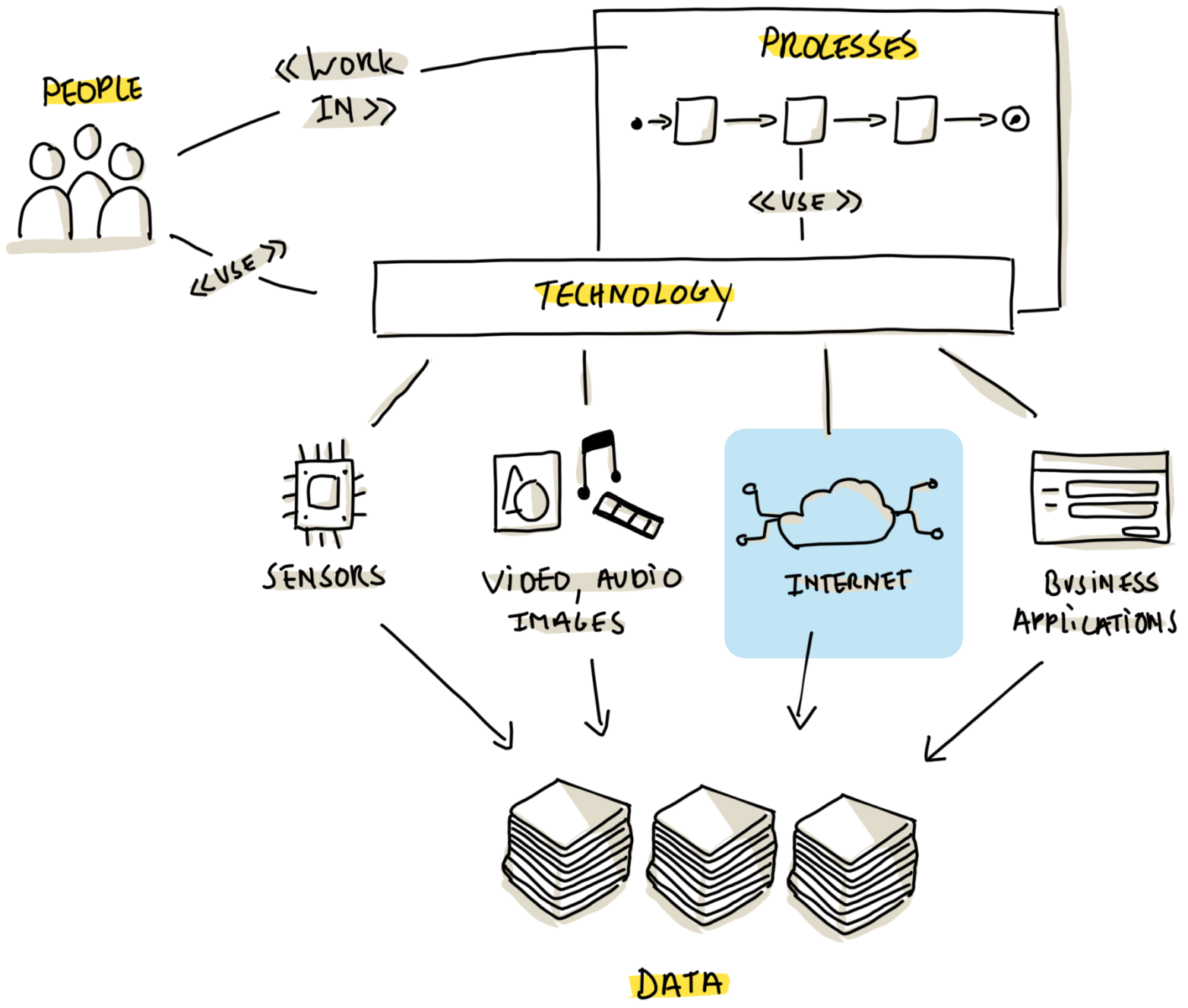
Document Scanning Apps





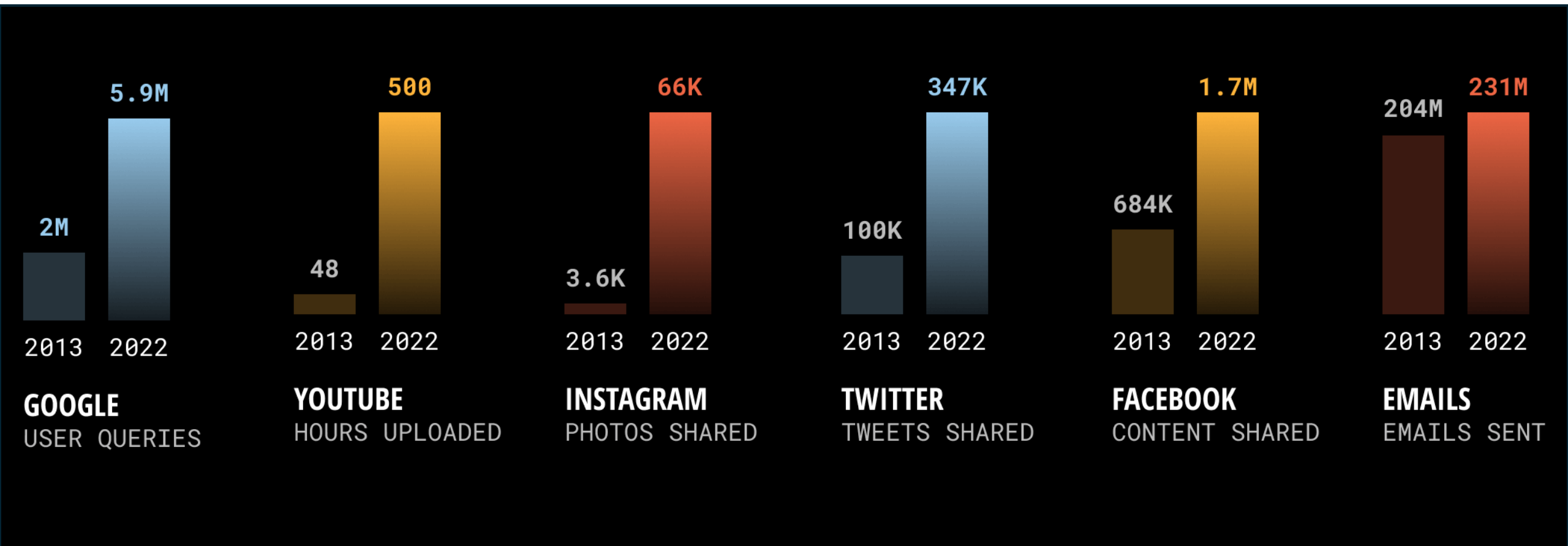
Digital Files







Big Data Producers: The Internet





1 MINUTE ON THE INTERNET



347,222
Stories uploaded



500 Hrs
Video Uploaded



41.7 Mn
Messages



319
New Users



4 Million
Searches



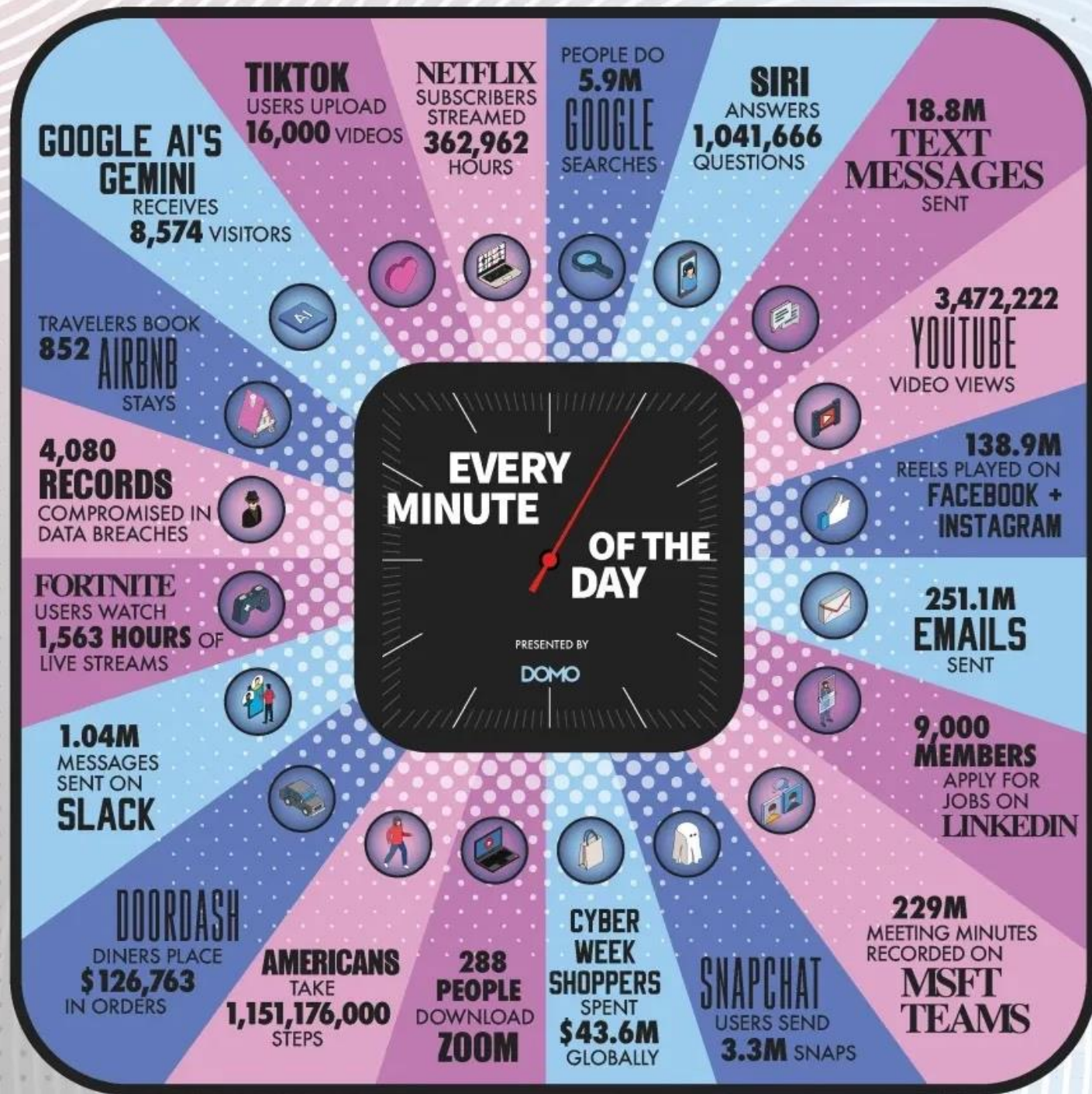
28 tracks
added to library

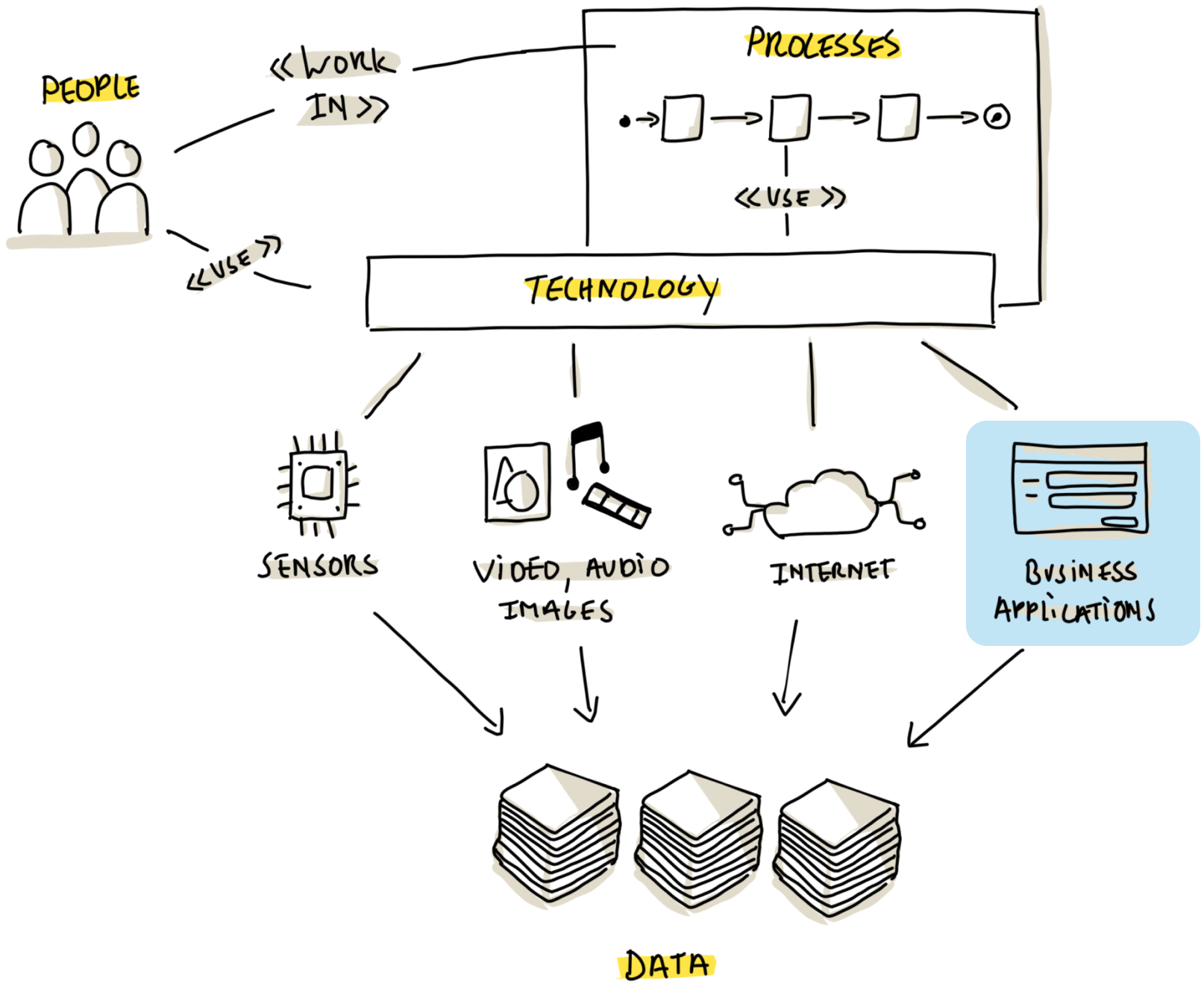


6660
Packages Shipped



404,444 Hrs
Streamed







Business Applications

- Business data is collected by **Business Applications**
- Examples:
 - **CRM systems**
 - **ERP systems**
 - **Service Desk Management system**
 - ...



My Open Opportunities | Salesforce

Secure https://acmecorp.lightning.force.com/one/one.app#/sObject/Opportunity/list?filterName=00BB0000002013TMAQ

Search Opportunities and more...

Lightning Sales Co... Opportunities Bennie Caudle Advanced Commun... Advanced Commun... Ben Seale

OPPORTUNITIES * My Open Opportunities

77 items • Sorted by Amount • Updated a minute ago

Qualification (37) Needs Analysis (30) Proposal/Quote (4) Negotiation (4) Closed Won (2)

Qualification (37)	Needs Analysis (30)	Proposal/Quote (4)	Negotiation (4)	Closed Won (2)
\$1,747,375	\$1,396,291	\$328,250	\$207,350	\$124,700
<p>Allied Technologies - Ne... \$20,000.00 Allied Technologies Qualification</p> <p>Amptech Corporation - A... \$61,550.00 Amptech Corporation Qualification</p> <p>Big Sky & Sons - Add-On ... \$72,000.00 Big Sky & Sons Qualification</p> <p>Biospan, LLC. - Services - ... \$77,700.00 Biospan, LLC. Qualification</p>	<p>Advanced Communicati... \$40,000.00 Advanced Communications Needs Analysis</p> <p>Allied Technologies - Ad... \$26,000.00 Allied Technologies Needs Analysis</p> <p>Amptech Corporation - A... \$81,100.00 Amptech Corporation Needs Analysis</p> <p>Amptech Corporation - ... \$41,000.00 Amptech Corporation Needs Analysis</p>	<p>Employnet - Services - 1... \$100,500.00 Employnet Proposal/Quote</p> <p>Haven Enterprises - Add-... \$61,500.00 Haven Enterprises Proposal/Quote</p> <p>Southern Sound Co. - Ad... \$124,600.00 Southern Sound Co. Proposal/Quote</p> <p>Towson Inc. - Add-On Bu... \$41,650.00 Towson Inc. Proposal/Quote</p>	<p>Opportunity Resources I... \$41,000.00 Opportunity Resources Inc Negotiation</p> <p>Tyconet - Add-On Busine... \$70,000.00 Tyconet Negotiation</p> <p>Opportunity Resources I... \$55,250.00 Opportunity Resources Inc Negotiation</p> <p>Tyconet - New Business - ... \$41,100.00 Tyconet Negotiation</p>	<p>Big Sky & Sons - Add-On ... \$66,000.00 Big Sky & Sons Closed Won</p> <p>Tech Labs - Add-On Busi... \$58,700.00 Tech Labs Closed Won</p>

Charts

Pipeline By Stage

Qualification	1,963,125
Needs Analysis	1,585,019.5
Value Proposition	1,341,890.4

Phone History Notes

[Source: [Salesforce Lightning](#)]



ERP: Purchasing Example in ODOO

Purchase Orders Products Reporting Configuration

New Requests for Quotation 1-24 / 24

All RFQs	10 To Send	3 Waiting	8 Late	Avg Order Value	\$ 3,529.01	Purchased Last 7 Days	\$ 24,916.48
My RFQs	6	3	5	Lead Time to Purchase	3.25 Days	RFQs Sent Last 7 Days	3

<input type="checkbox"/>	Reference	Vendor	Buyer	Order Deadline	Activities	Source Docu...	Total	Status	
<input type="checkbox"/>	<input type="checkbox"/> ☆ P00024	Azure Interior		In 4 days	Call	OP/00003	\$ 316.25	RFQ	
<input type="checkbox"/>	<input type="checkbox"/> ☆ P00023	Ready Mat	Mitchell Admin	Today	Confirm order		\$ 495.00	RFQ Sent	
<input type="checkbox"/>	<input type="checkbox"/> ☆ P00022	Gemini Furniture	Mitchell Admin	Today	Reminder		\$ 460.00	RFQ Sent	
<input type="checkbox"/>	<input type="checkbox"/> ☆ P00021	Azure Interior	Mitchell Admin	Today	Email		\$ 420.00	RFQ Sent	
<input type="checkbox"/>	<input type="checkbox"/> ☆ P00020	Gemini Furniture				Replenishment R...	\$ 230.58	Purchase Order	
<input type="checkbox"/>	<input type="checkbox"/> ☆ P00019	The Jackson Group				Replenishment R...	\$ 1,667.50	Purchase Order	
<input type="checkbox"/>	<input type="checkbox"/> ☆ P00018	Wood Corner			Send shipping...		\$ 46.00	Purchase Order	
<input type="checkbox"/>	<input type="checkbox"/> ☆ P00017	Wood Corner					\$ 276.00	Purchase Order	
<input type="checkbox"/>	<input type="checkbox"/> ☆ P00016	YourCompany, Jo...	Marc Demo				\$ 57.50	Purchase Order	
<input type="checkbox"/>	<input type="checkbox"/> ☆ P00015	Ready Mat	Mitchell Admin				\$ 6,596.40	Purchase Order	



Service Desk Management: Example

Jira Service Management Your work **Projects** Filters Dashboards People Plans Apps **Create** 9+

ITSM Service Desk
Service project

- Queues**
- Service requests
- Incidents
- Problems
- Changes

OPERATIONS

- Services
- Alerts
- On-call

KNOWLEDGE & INSIGHTS

- Knowledge base
- Reports

CHANNELS & PEOPLE

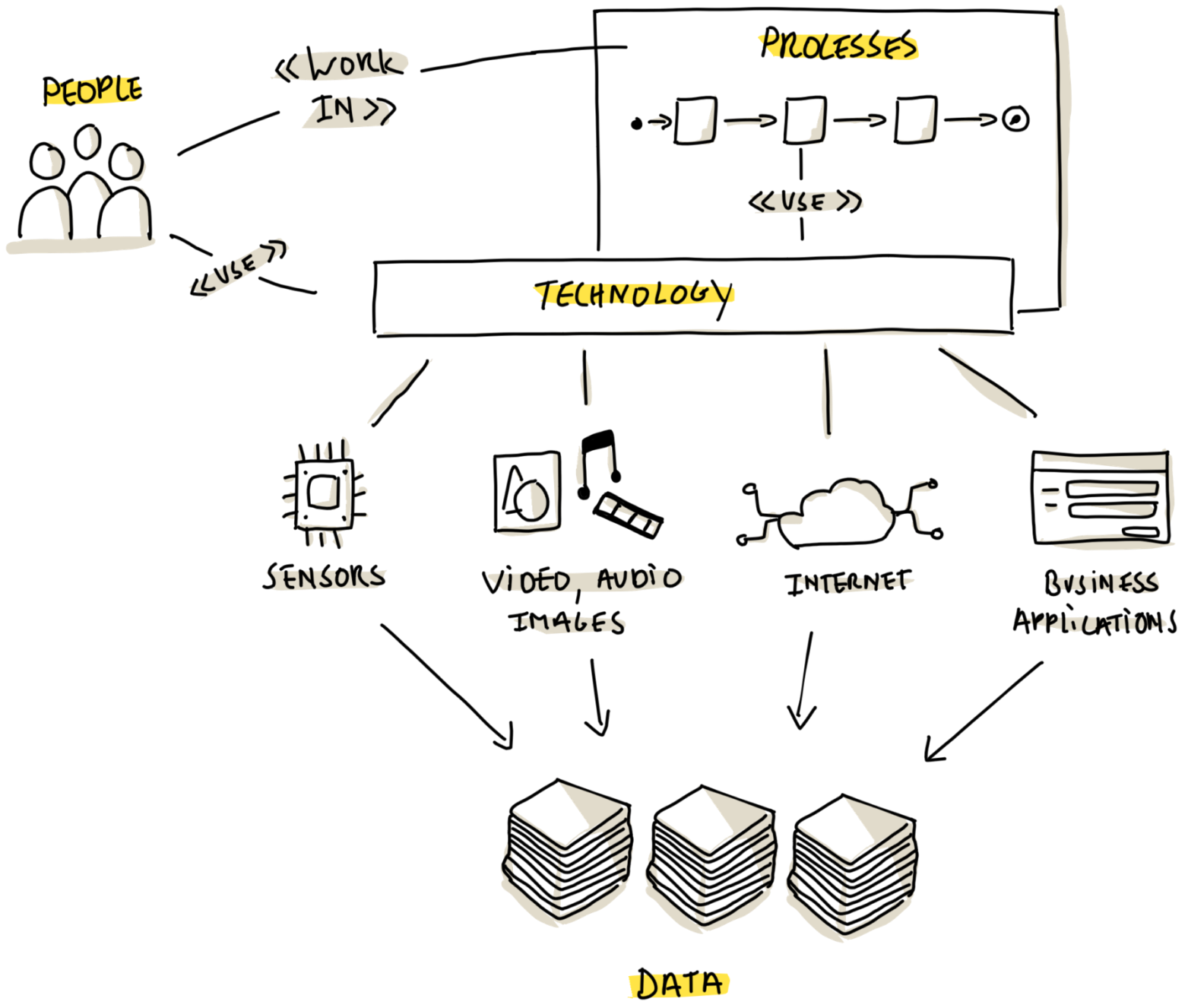
- Ticket channels
- Customers
- Invite team

Project settings

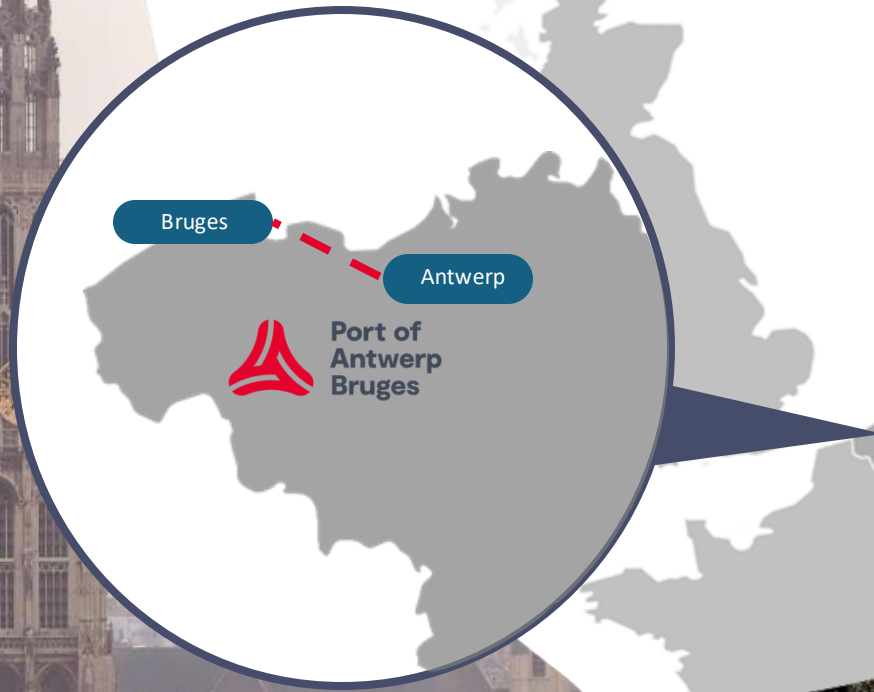
Projects / ITSM Service Desk / All tickets

All open tickets

<input type="checkbox"/>	Request Type	Key	Summary	Reporter	Assignee	Status	Created	Time to resolu... ↑	P
<input type="checkbox"/>	Request admin access	ITSM-1324	Admin access to Jira	Polly ProductManager	Sammy Servi...	WAITING FOR SUPPORT	24/Sep/20	-15m ⌚	↑
<input type="checkbox"/>	Report a system problem	ITSM-1342	Banc.ly Inc is slow	Serena ServiceDeskManager	Sammy Servi...	WORK IN PROGRESS	25/Sep/20	3h 44m ⌚	↑
<input type="checkbox"/>	Report a system problem	ITSM-1339	Can't access POS System	Sammy ServiceDeskAgent	Sammy Servi...	WORK IN PROGRESS	25/Sep/20	3h 44m ⌚	↑
<input type="checkbox"/>	Report broken hardware	ITSM-1333	Cant access webcam	Darrel DevManager	Sammy Servi...	WORK IN PROGRESS	25/Sep/20	3h 44m ⌚	↑
<input type="checkbox"/>	Report a system problem	ITSM-1331	Can't access Trello	Sammy ServiceDeskAgent	Serena Servi...	WORK IN PROGRESS	25/Sep/20	3h 44m ⌚	↑
<input type="checkbox"/>	Request admin access	ITSM-1338	Admin access to Jira	Polly ProductManager	Sammy Servi...	IN PROGRESS	25/Sep/20	7h 44m ⌚	↑
<input type="checkbox"/>	Get a guest wifi account	ITSM-1337	Guest access	Darla DevDirector	Sammy Servi...	IN PROGRESS	25/Sep/20	7h 44m ⌚	↑
<input type="checkbox"/>	Request new software	ITSM-1336	Add Office to Mac	Sammy ServiceDeskAgent	Serena Servi...	IN PROGRESS	25/Sep/20	7h 44m ⌚	↑
<input type="checkbox"/>	Request new hardware	ITSM-1335	Need new keyboard	Sandeep ServiceOwner	Serena Servi...	IN PROGRESS	25/Sep/20	7h 44m ⌚	↑
<input type="checkbox"/>	Set up VPN to the office	ITSM-1334	VPN Access	Christy ChangeManager	Serena Servi...	IN PROGRESS	25/Sep/20	7h 44m ⌚	↑
<input type="checkbox"/>	New mobile device	ITSM-1332	Need a new iPhone	Dante Developer	Sammy Servi...	WAITING FOR APPROVAL	25/Sep/20	7h 44m ⌚	↑
<input type="checkbox"/>	Get IT help	ITSM-1330	Help setting up my VPN	Carly ChiefExec	Sammy Servi...	IN PROGRESS	25/Sep/20	7h 44m ⌚	↑



A global port in the heart of Europe



One port
Two sites





2nd largest port in **Europe**



Number one **export** port in Europe

147 mio tons/year



Largest **car handling** port in Europe

3,160,000 cars/year



21,400 **Seagoing ships**/year



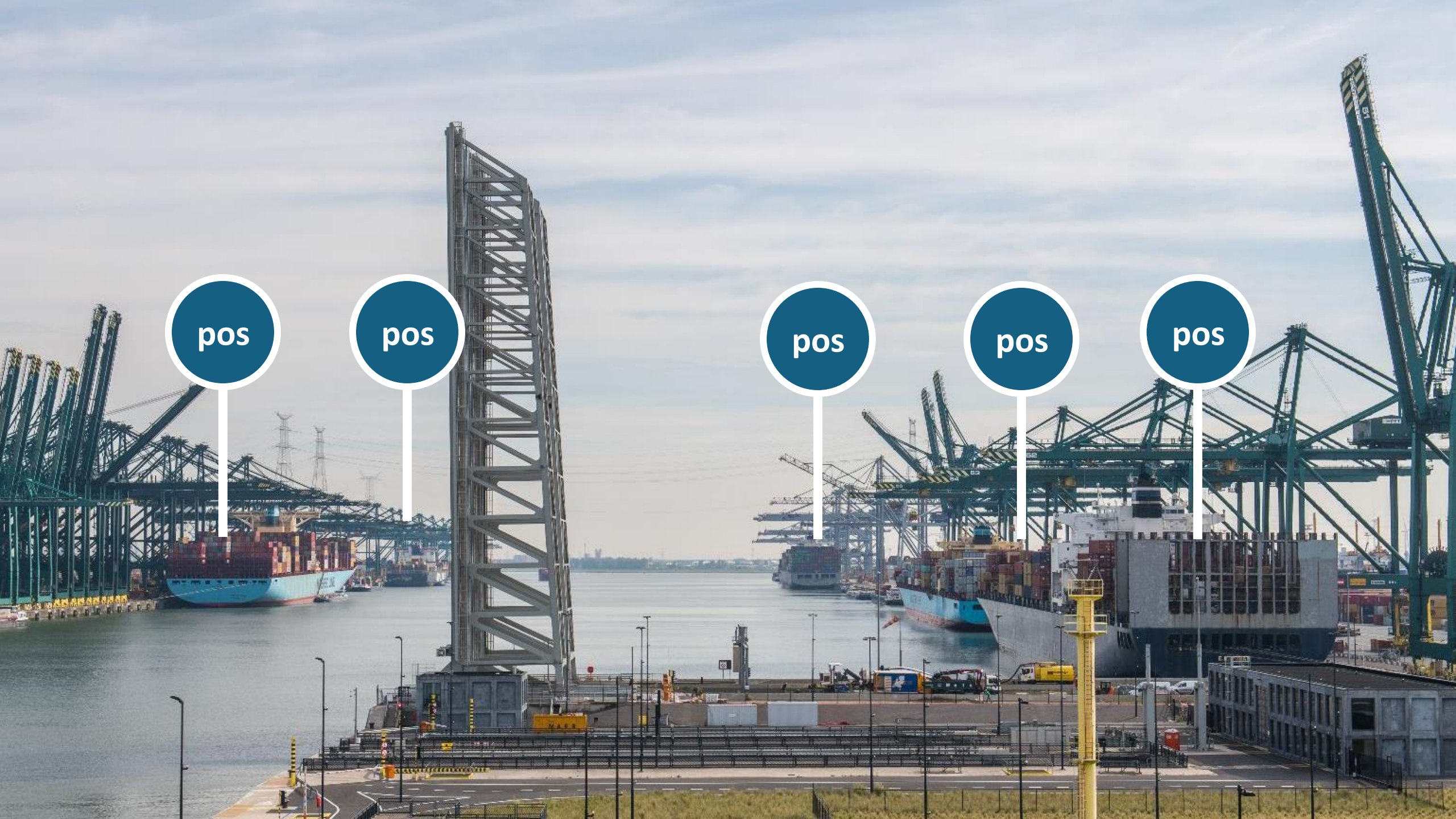
15% of EU **gas** market

Which Data is Available at Port Of Antwerp Bruges?



Data Is Everywhere at The Port





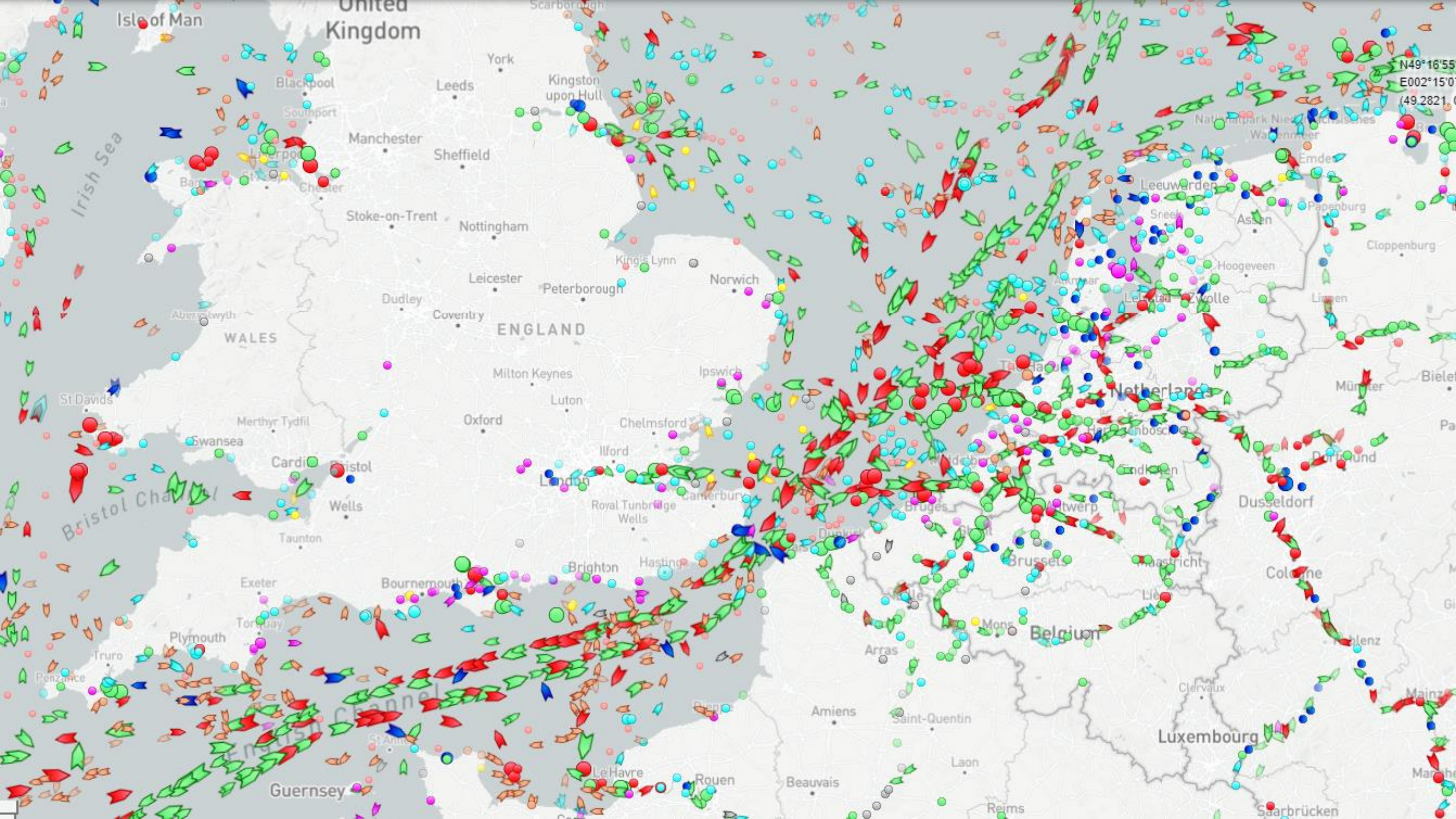
pos

pos

pos

pos

pos





Automatic Identification System (AIS)

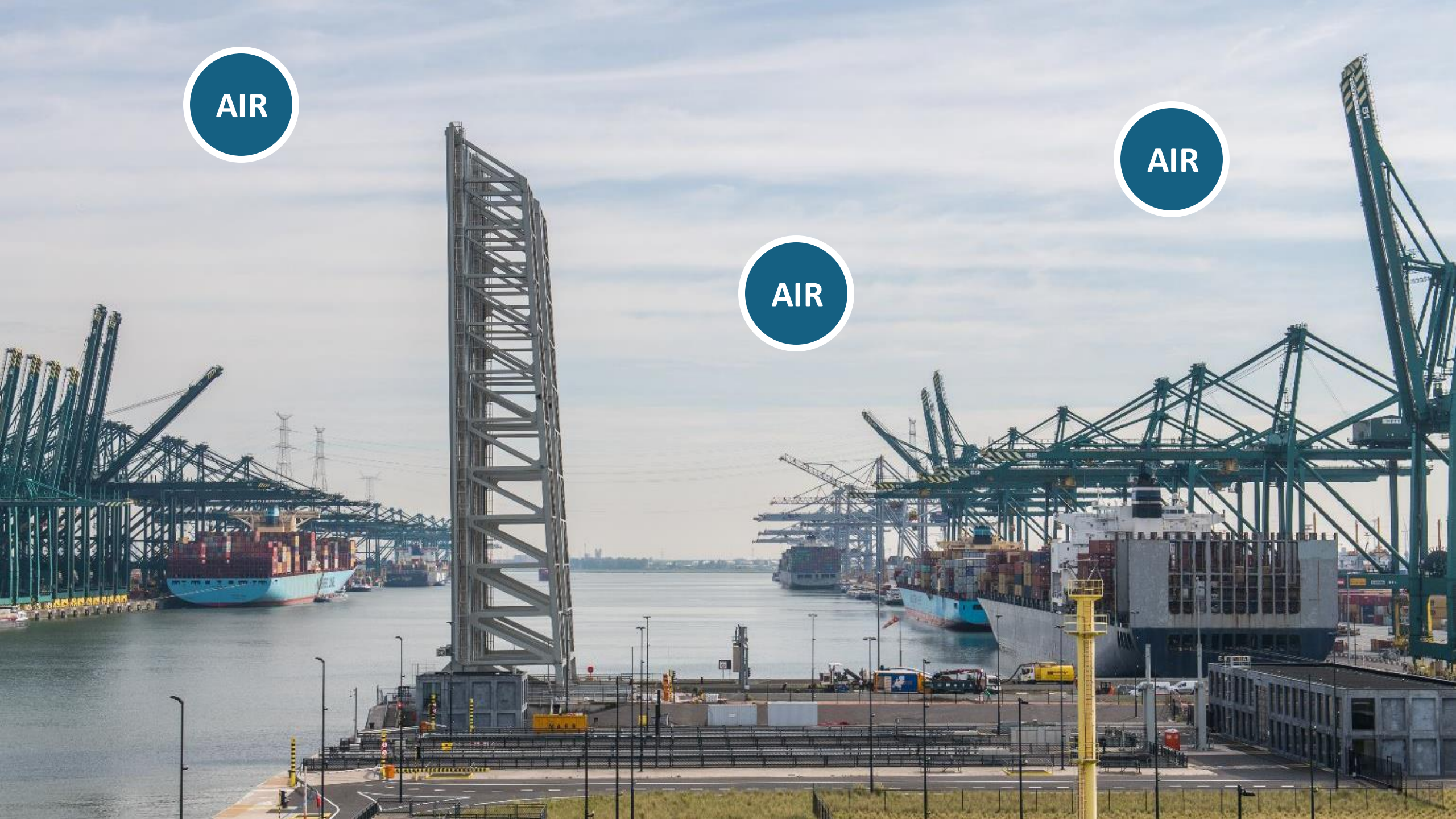
“Public” Data



AIR

AIR

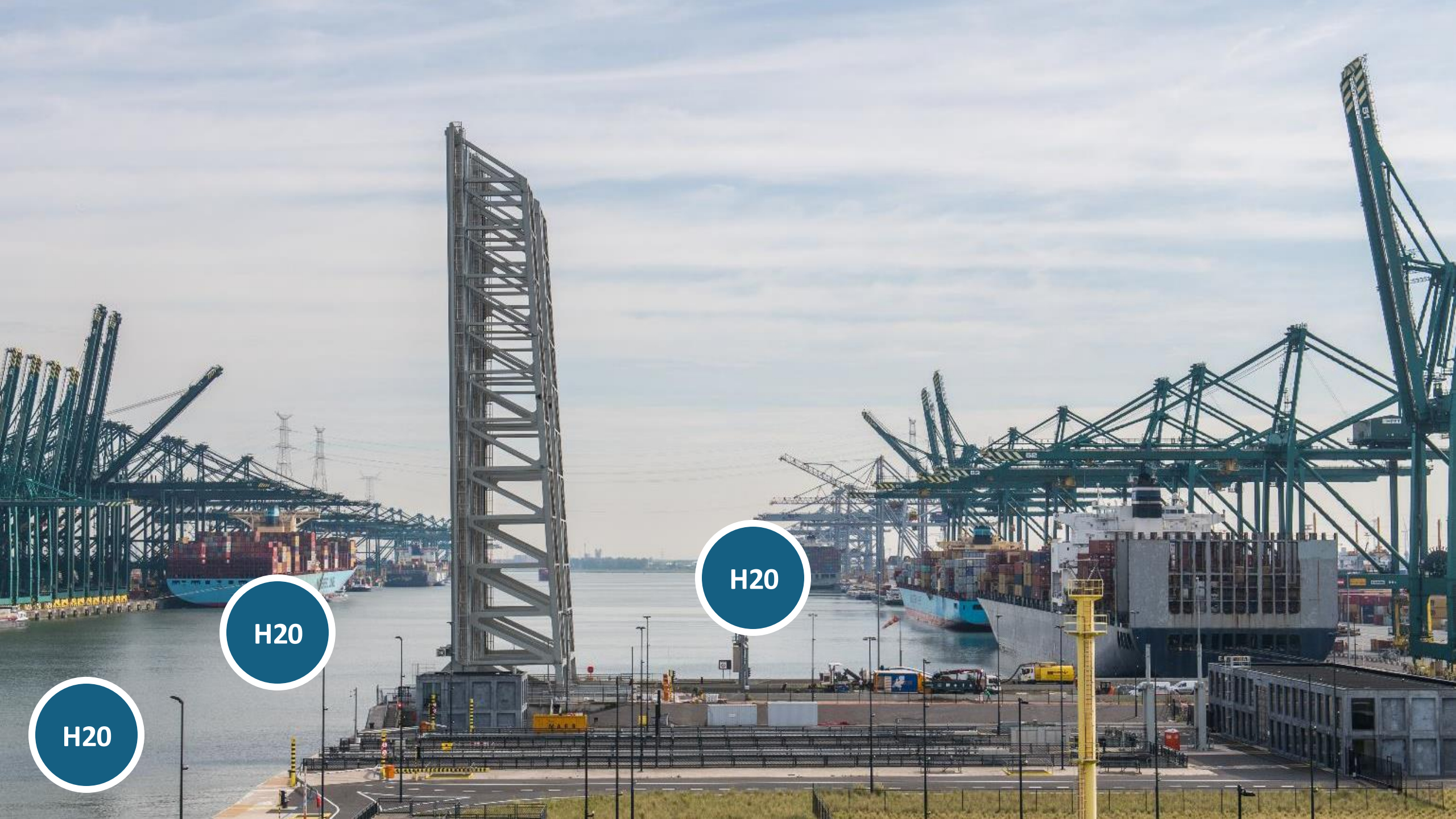
AIR





Air Quality Sensors





H2O

H2O

H2O



Water Level Sensors



Floating Waste



Port of
Antwerp
Bruges



BRIDGE

Bridges



Port of
Antwerp
Bruges

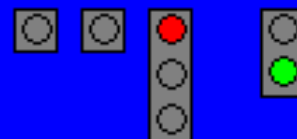
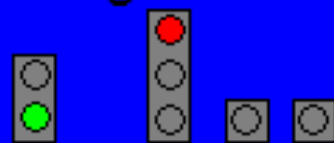


Westkant

● COM ERR WAGO PEB



Lillobrug



00:09:29



Oostkant

● COM ERR NOS LILLOBRUG

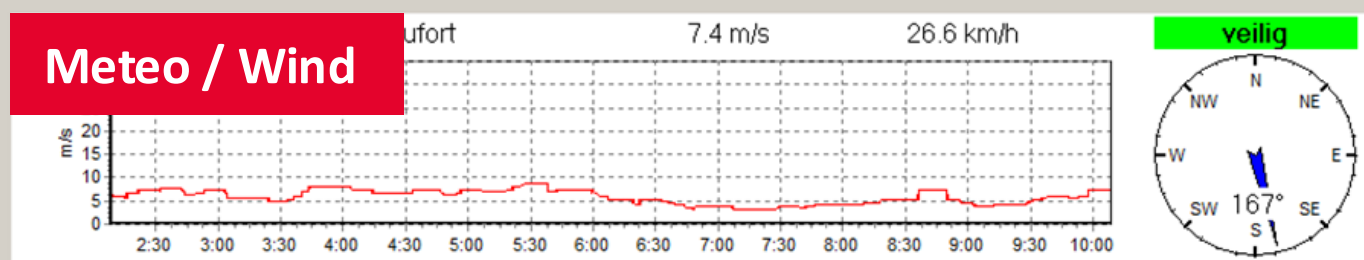
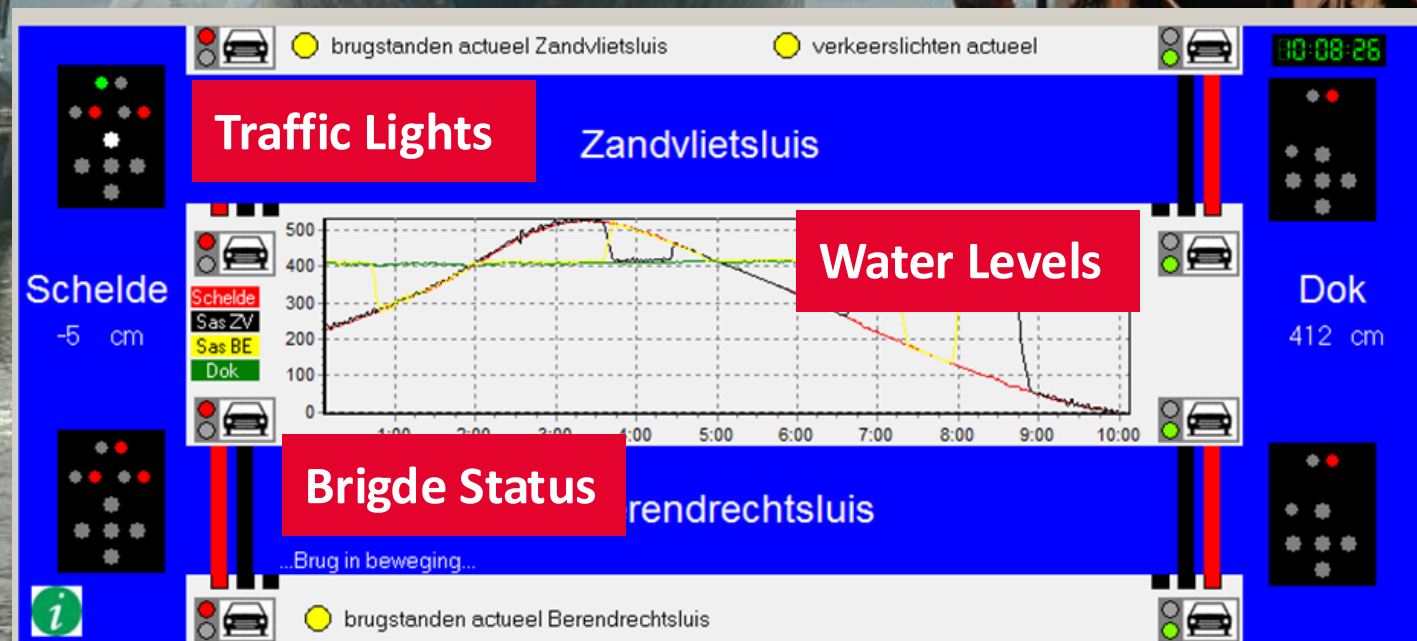


Locks



Port of
Antwerp
Bruges

Locks



ENERGY



Energy



Port of
Antwerp
Bruges

fluvius.

Data Is Everywhere at The Port



AIR

ENERGY

pos

pos

AIR

AIR

pos

pos

pos

BRIDGE

H2O

H2O

H2O

Data Is Everywhere at The Port

SENSOR DATA

DATA THAT WE OWN

DATA THAT WE USE

TRAFFIC
WASTE
VIDEO
LOCKS
BRIDGES
WATER

METEO
AIS ENERGY
AIR QUALITY

MANY MORE

APPLICATION DATA

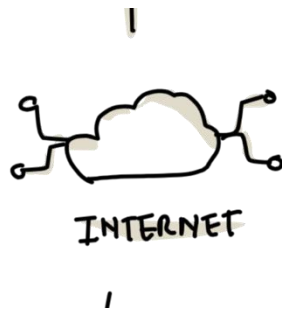
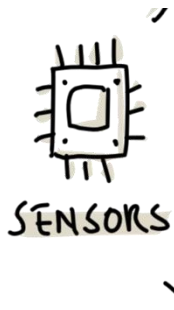
ERP CRM HR

ASSET MANAGEMENT

MANY MORE

OEFENING: Welke Data Providers leveren data voor jullie use case?

- Welke data providers?
- Welk type providers?



- Bespreek kort in groep – delen achteraf

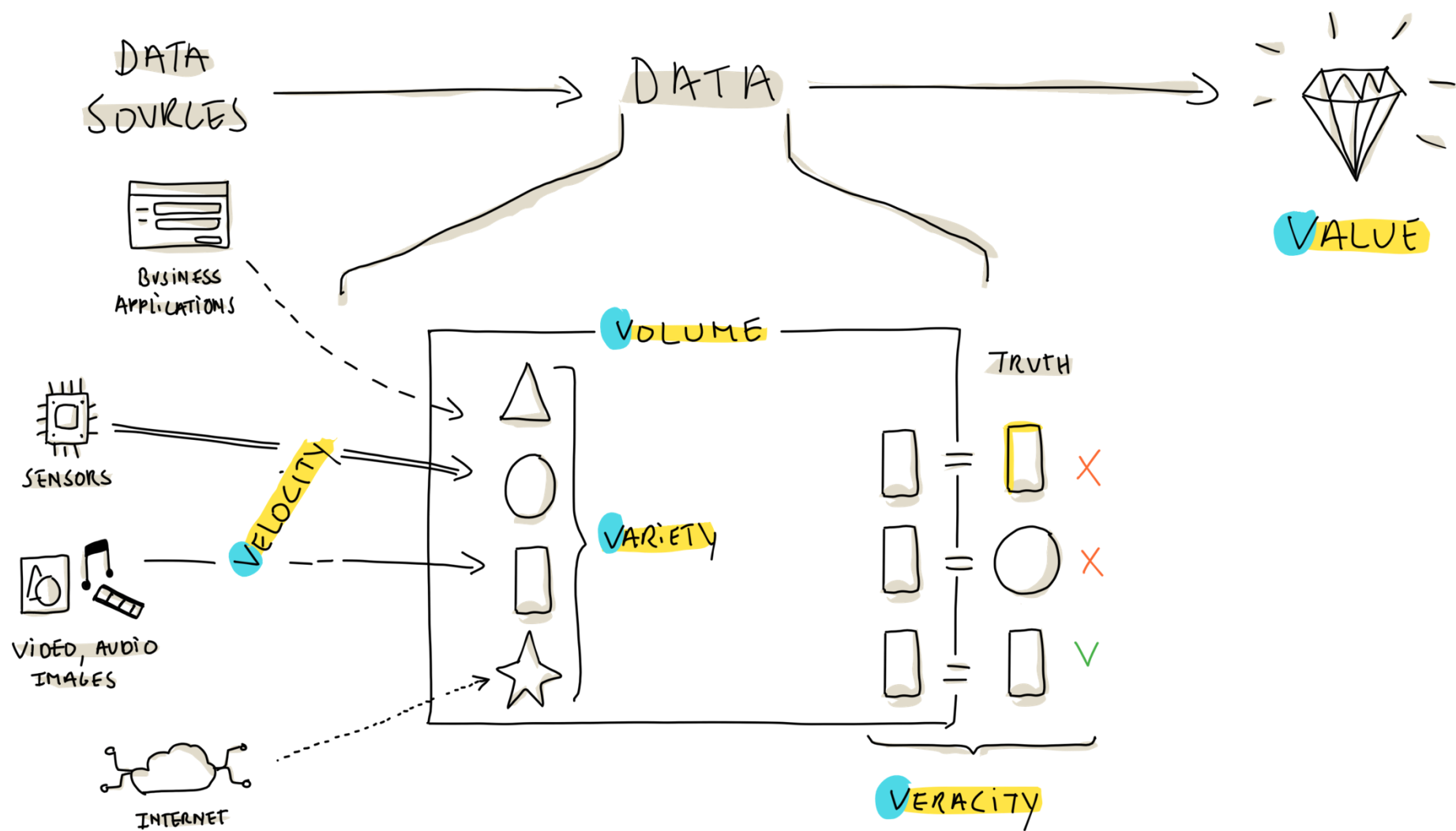


Data Management

- Defining Data
- Data Producers
- **Big Data Vs**
- The Bigger Picture
- Data Management
- Data Technology



5 Data V's



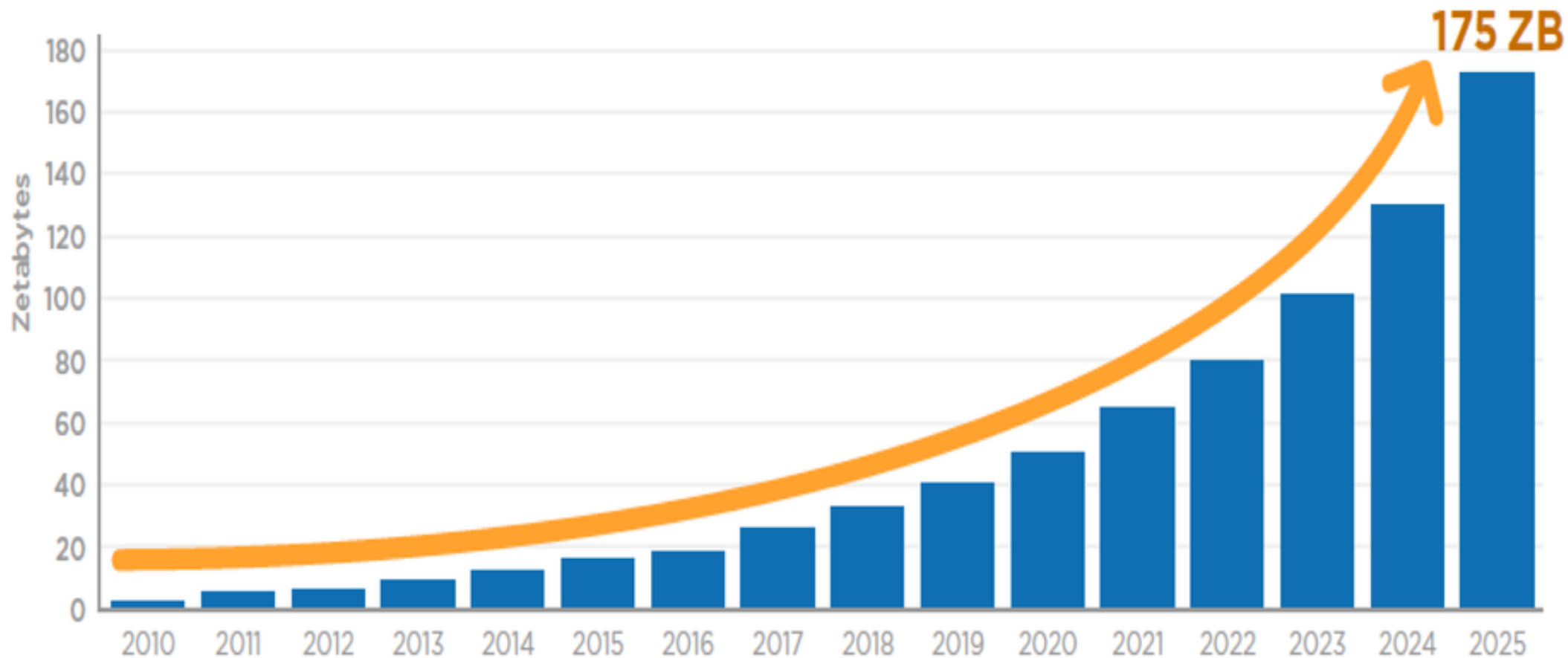


Volume - the size of Big Data

- Big data analytics entails handling and analyzing **vast amounts of data**.
- To effectively work with such massive datasets, **specialized tools and infrastructure are necessary** for capturing, storing, managing, cleaning, transforming, analyzing, and reporting the data.



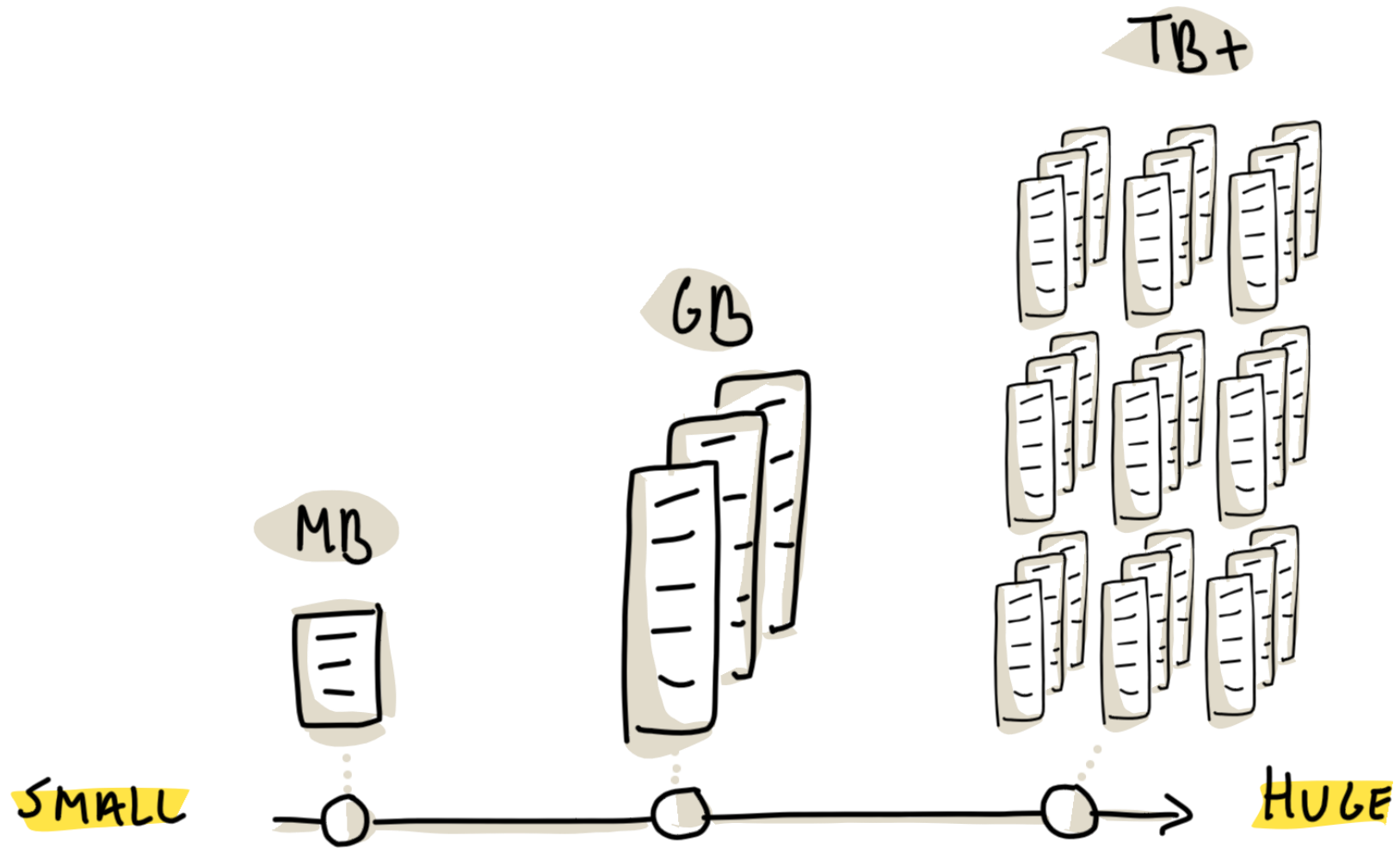
Volume



Source: Data Age 2025, sponsored by Seagate with data from IDC Global DataSphere, Nov 2018

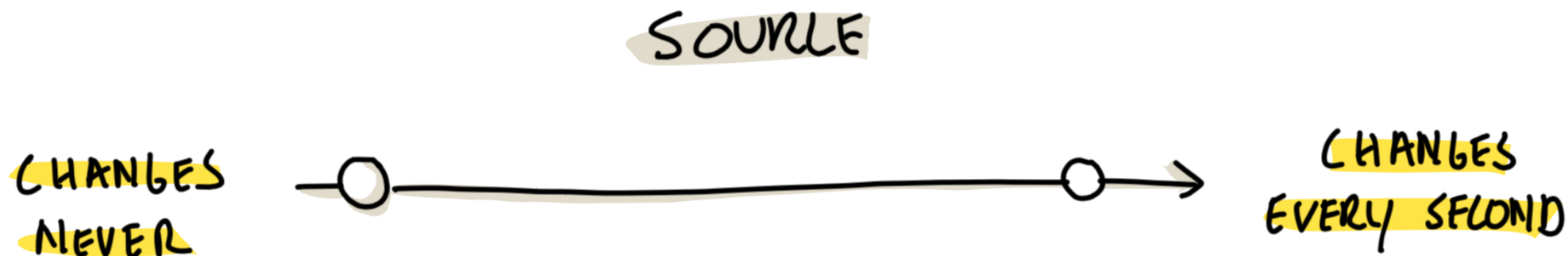


Volume



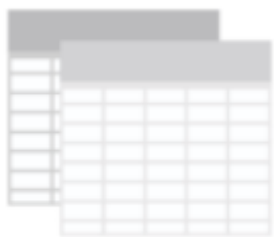
Velocity - the speed at which data is growing

- Velocity denotes the speed at which data is generated.
- To keep up with the rapid generation of data, systems for processing and analyzing data must possess **sufficient capacity to handle the influx of data and deliver timely, actionable insights.**





Variety - the different structures of data



Structured
Table Data

STRUCTURED



XML



JSON



SENSOR

SEMI STRUCTURED



Text



Image



Video



Audio

UNSTRUCTURED



Variety – Structured Data

- Constituted of elements that are **readily addressable to facilitate effective analysis**.
- This category typically arranges **data into tabular structures** with rows and columns.
- Structured data presents a straightforward starting point for data analysis.
- Examples encompass relational data, CSV files, and spreadsheets.

Tabular Data

ID	TOTAL ACTIONS	ACTION 1	ACTION 2	TOTAL TIME
10	120	80	40	0:50:05
11	255	130	125	1:40:03
12	180	100	80	1:20:19
13	305	205	100	1:58:58
14	71	50	21	0:35:41
15	418	310	108	2:08:18
16	222	150	72	1:32:58



Variety – Semi-structured Data

- This data type **lacks the rigid tabular format**
- Retains some inherent organizational traits that render it **more analyzable compared to unstructured data.**
- Illustrations include XML data and JSON data.



Variety – Semi-structured Data

XML

vs.

JSON

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <endereco>
3   <cep>31270901</cep>
4   <city>Belo Horizonte</city>
5   <neighborhood>Pampulha</neighborhood>
6   <service>correios</service>
7   <state>MG</state>
8   <street>Av. Presidente Antônio Carlos, 6627</street>
9 </endereco>
```

```
1 {
2   "endereco": {
3     "cep": "31270901",
4     "city": "Belo Horizonte",
5     "neighborhood": "Pampulha",
6     "service": "correios",
7     "state": "MG",
8     "street": "Av. Presidente Antônio Carlos, 6627"
9   }
10 }
```

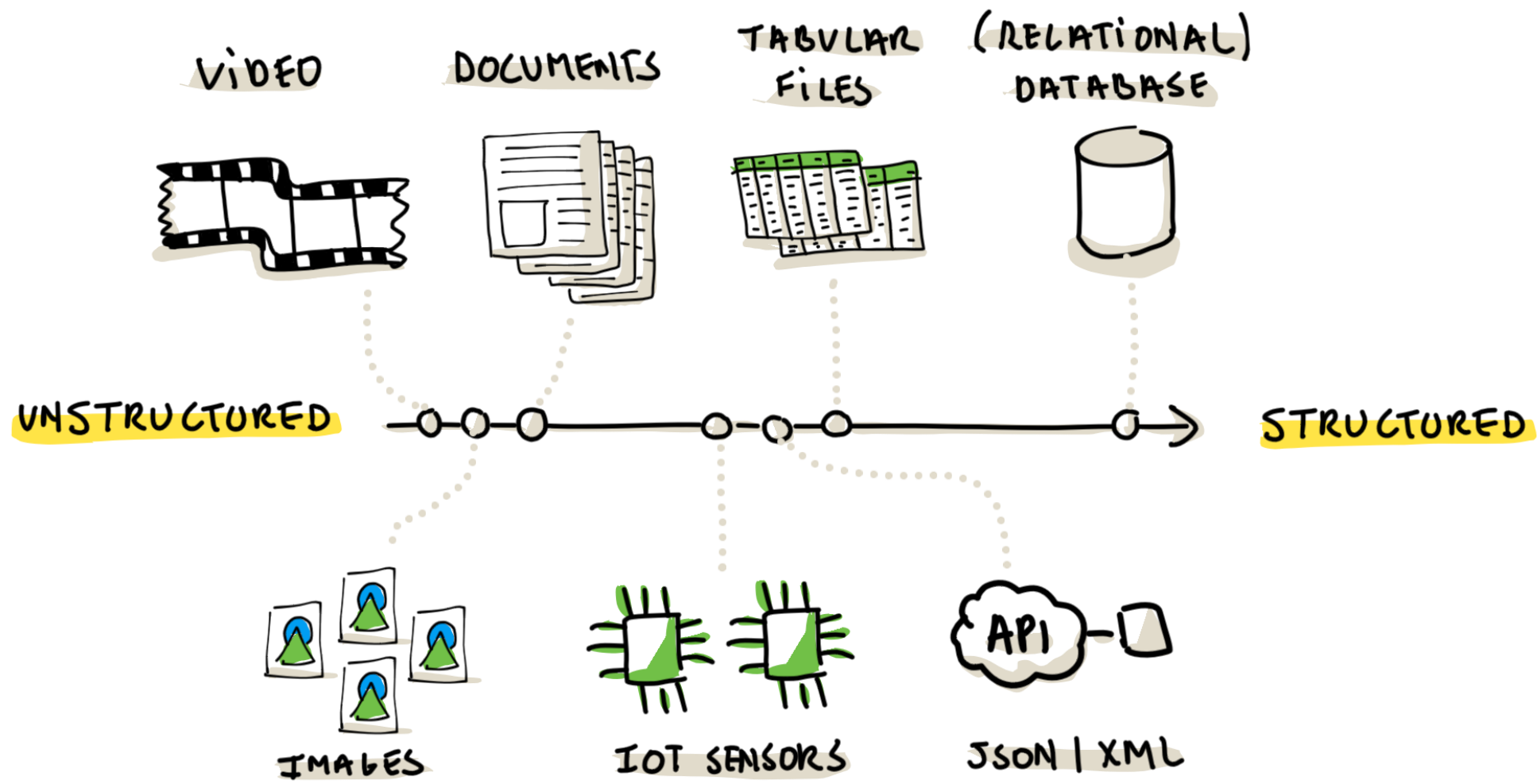


Variety – Unstructured Data

- Data lacking predefined organization or structure.
- The absence of clear formatting renders this data type intricate to analyze.
- Examples entail documents like Word files, PDFs, textual content, images, videos, and audio recordings.

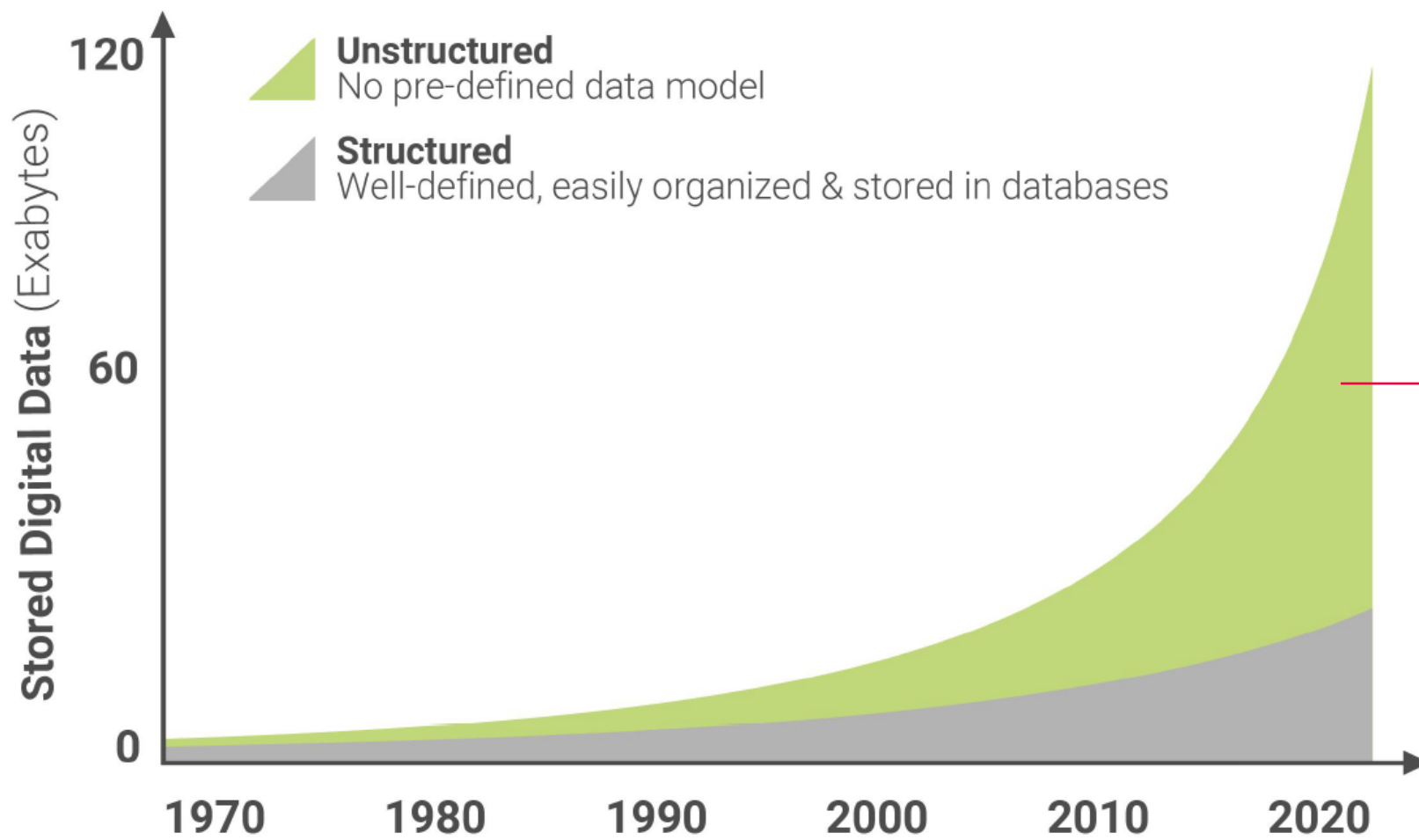


Variety





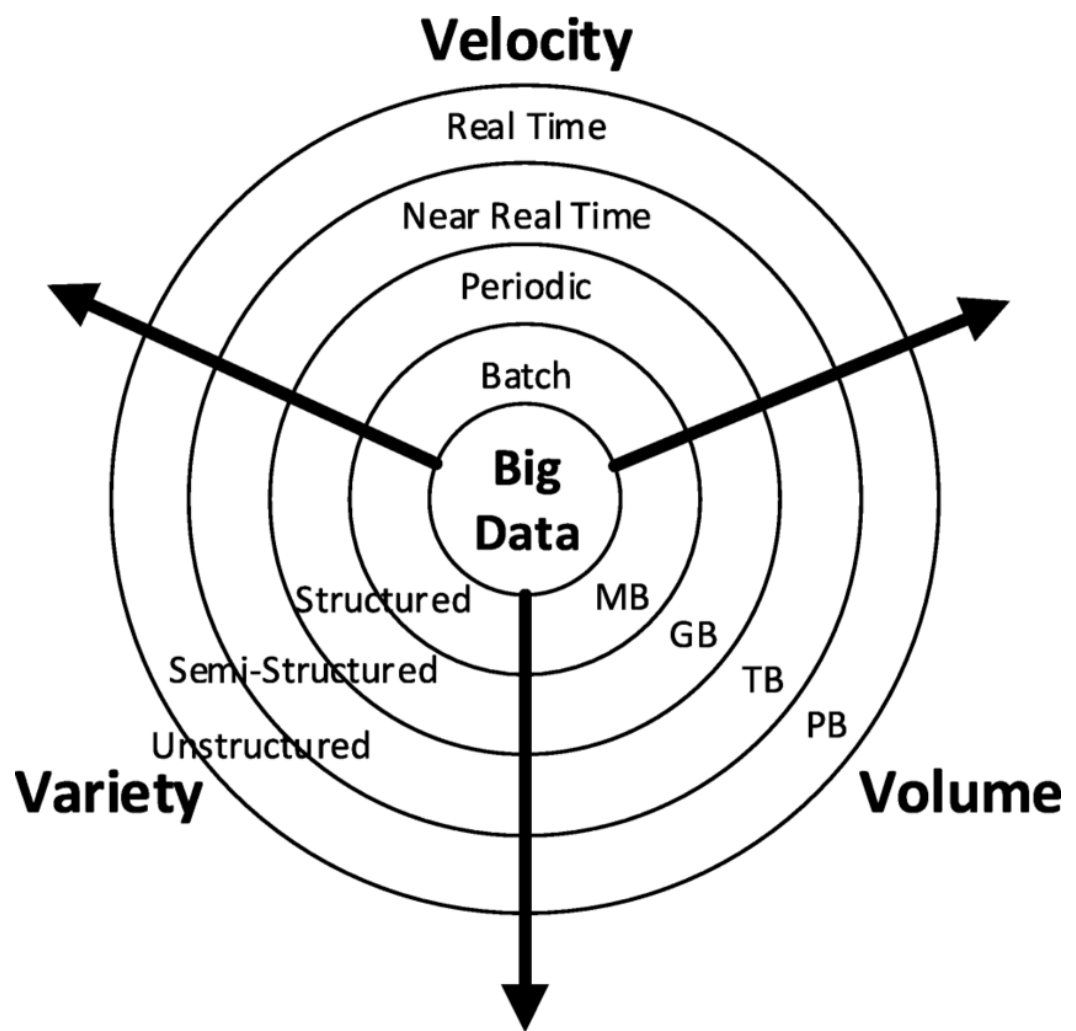
Variety / Velocity / Volume



Exponential growth of unstructured data!



Variety / Velocity / Volume





Veracity – Accuracy or truthfulness of data?

- Veracity pertains to the **accuracy and authenticity of the data**.
- Data must undergo **validation** to ensure that it accurately represents essential business functions.
- Any data manipulation does not compromise the data's accuracy.



Value – How useful is the data?





Value – How useful is the data?

- Big data must generate **value**.
- The insights derived from an analysis should for example provide meaningful guidance for **improving operations, enhancing customer service, or creating other forms of value**.

OEFENING: De 5 V's

- Wat zijn de belangrijkste eigenschappen (5 V's) van de data voor jullie use cases?
 1. Velocity
 2. Volume
 3. Variety
 4. Veracity
 5. Value
- Overleg kort in groep

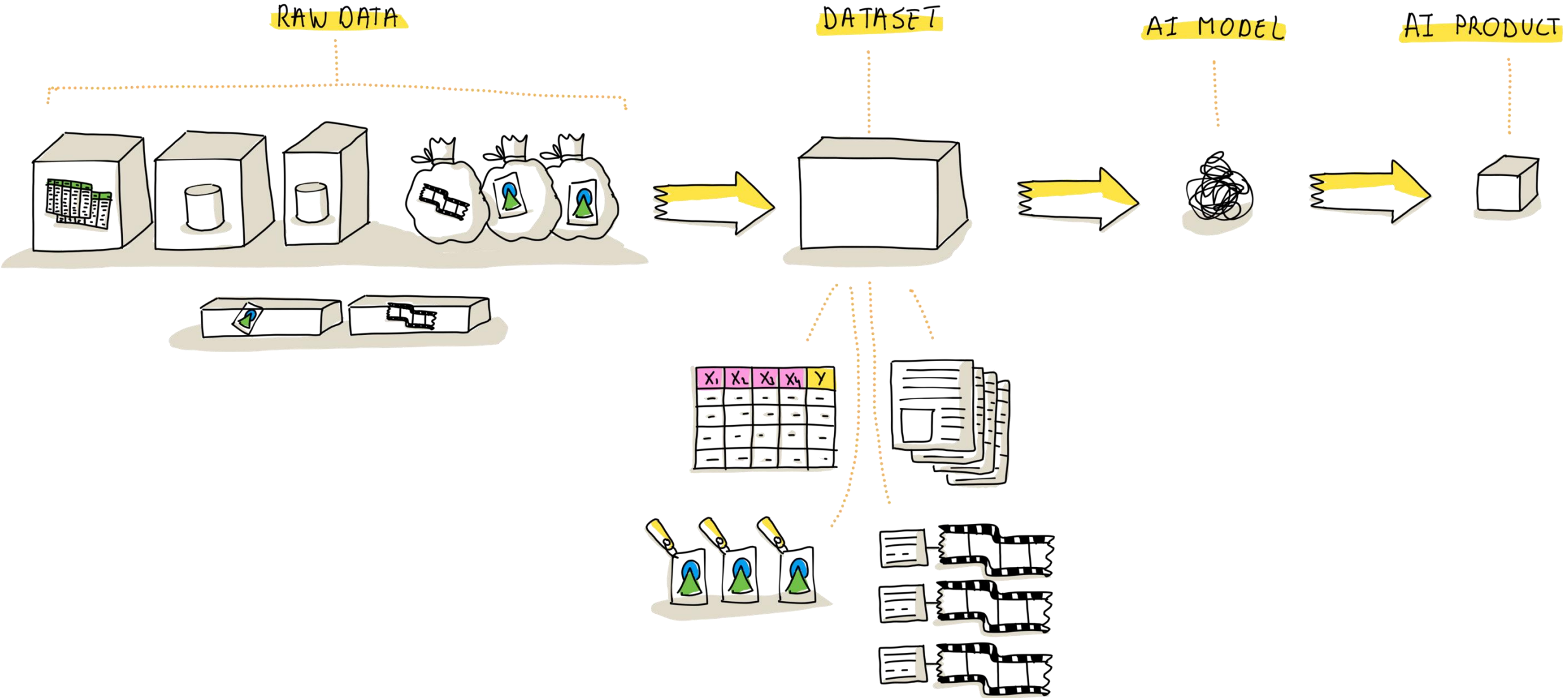


Data Management

- Defining Data
- Data Producers
- Big Data Vs
- **The Bigger Picture**
- Data Management
- Data Technology



The Bigger Picture





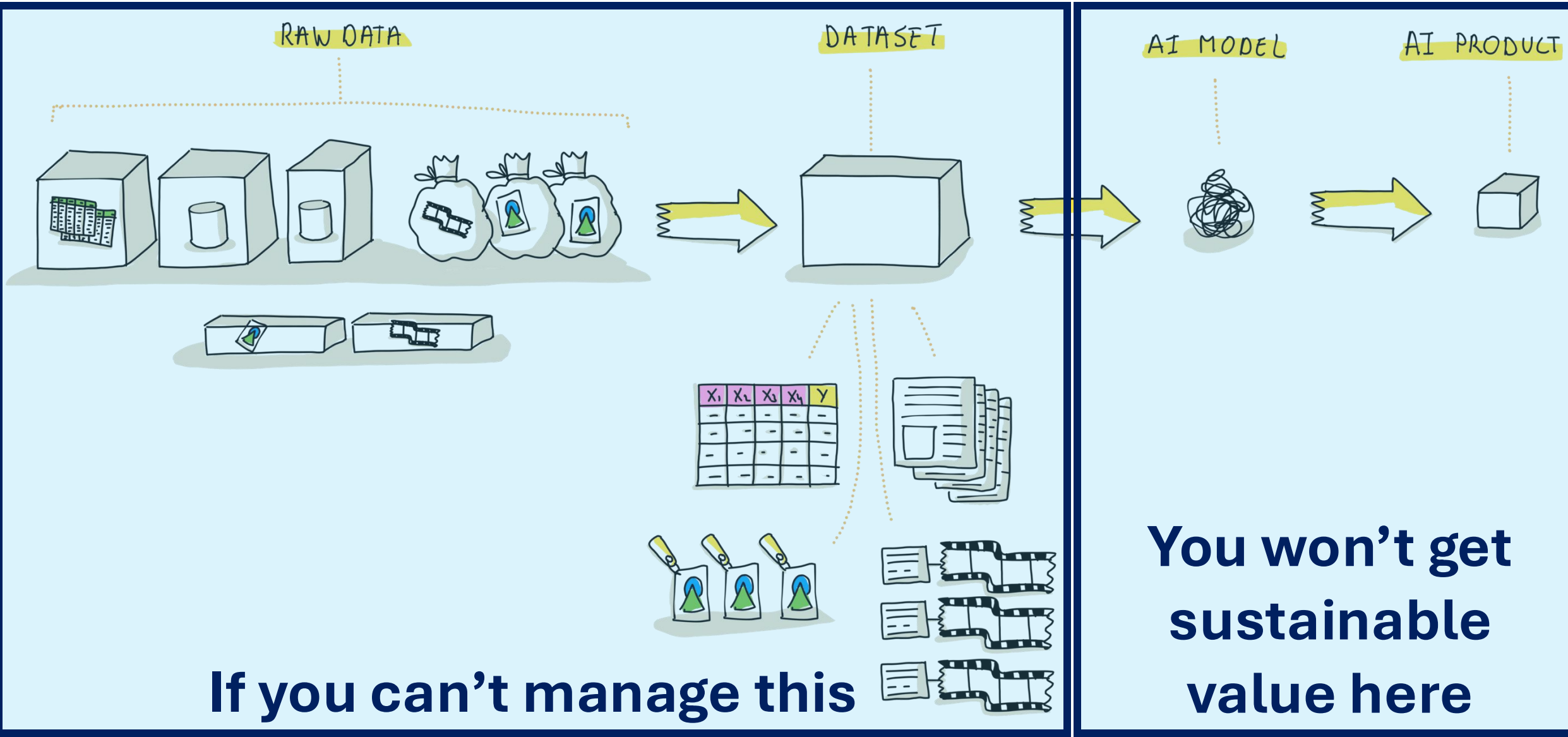
The Bigger Picture

RAW DATA

DATASET

AI MODEL

AI PRODUCT



If you can't manage this

You won't get sustainable value here



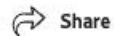
CIO JOURNAL

Rush to Use Generative AI Pushes Companies to Get Data in Order

Data management is under the spotlight again as companies seek to out-innovate competitors with large language models

By [Belle Lin](#) [Follow](#)

June 8, 2023 7:00 am ET



Share



Resize



Listen (2 min)



Corporate technology chiefs are under pressure to ensure companies' data is stored, filtered, and protected for use with AI. PHOTO: I-HWA CHENG/BLOOMBERG NEWS



Barriers to Preventing Organizations From Delivering More Value From Data, Analytics, and AI

Q11: Please rank all barriers preventing you from delivering more value from data, analytics, and AI, where 1 is the biggest/most present barrier (Note: Other (n=41) not shown):

% Ranking Each 'Rank 1 - Biggest/Most Present Barrier,' 'Rank 2,' or 'Rank 3'

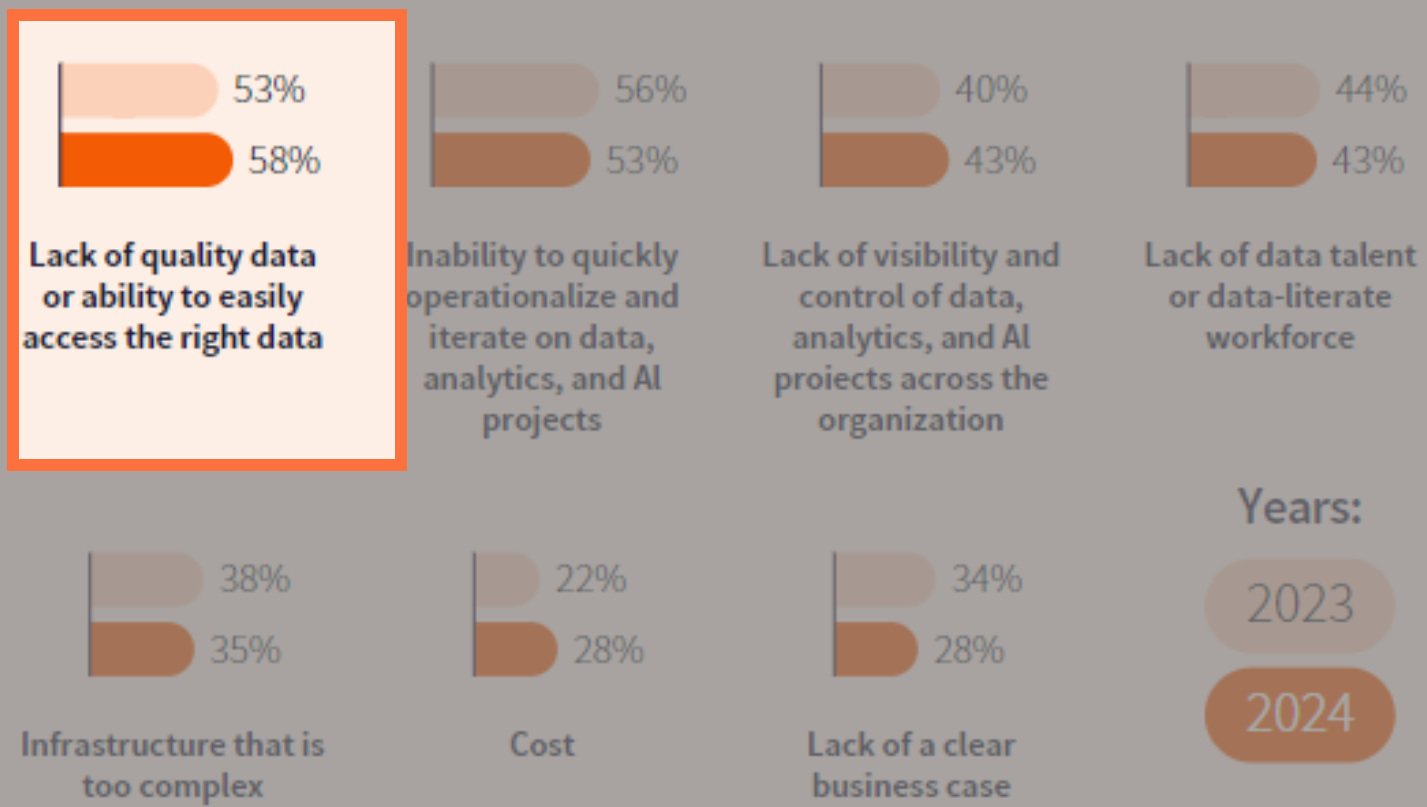




Barriers to Preventing Organizations From Delivering More Value From Data, Analytics, and AI

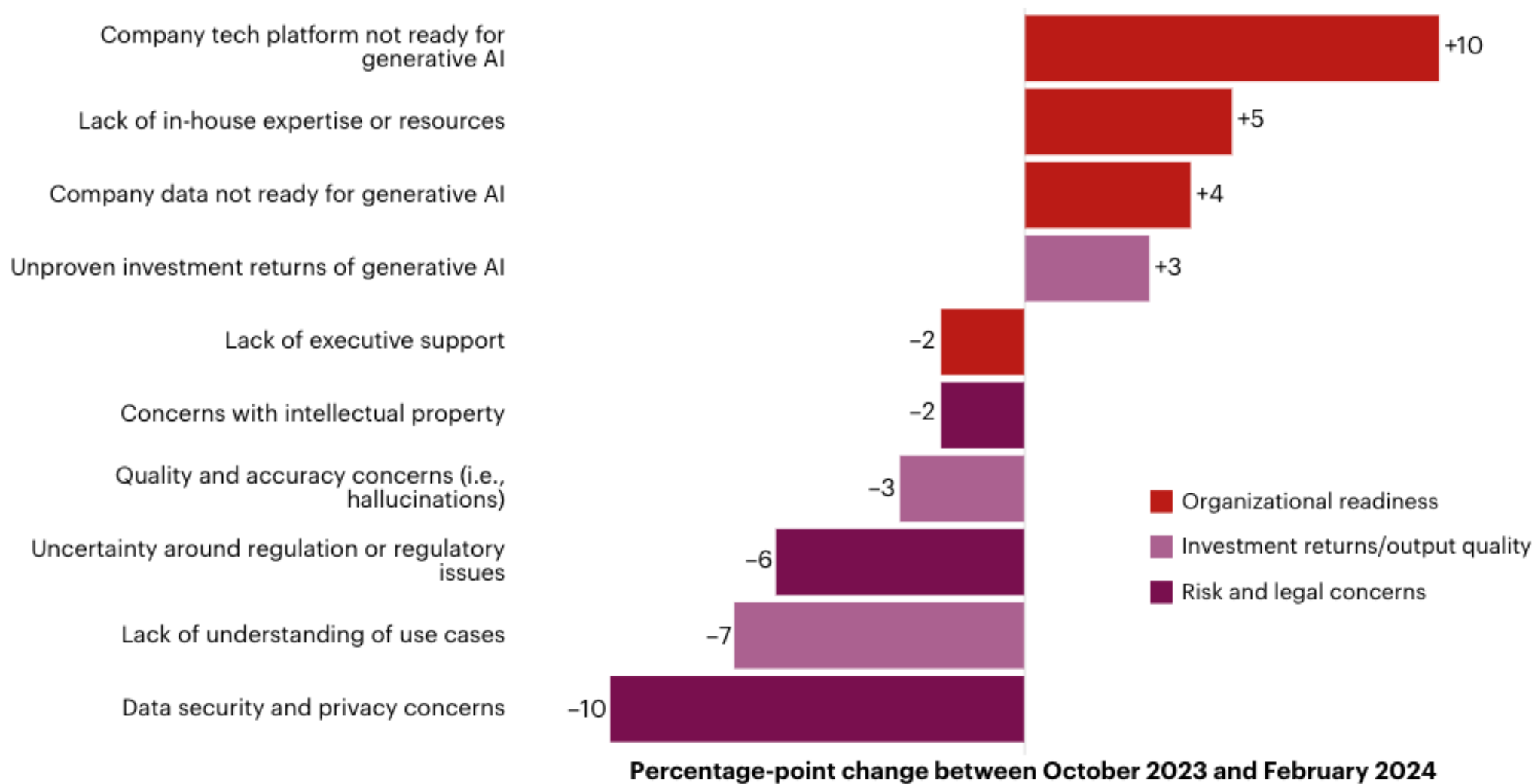
Q11: Please rank all barriers preventing you from delivering more value from data, analytics, and AI, where 1 is the biggest/most present barrier (Note: Other (n=41) not shown):

% Ranking Each 'Rank 1 - Biggest/Most Present Barrier,' 'Rank 2,' or 'Rank 3'





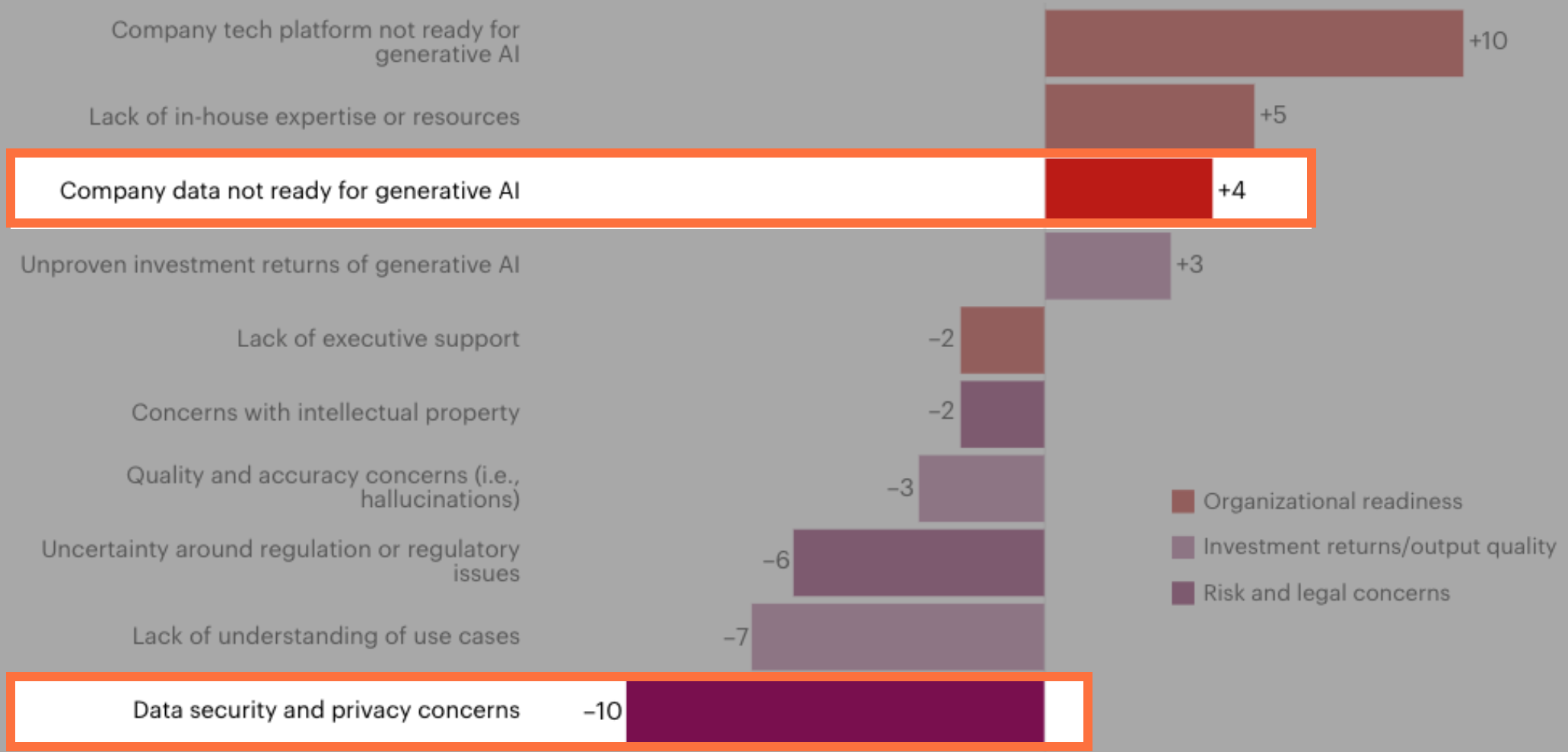
What are the top reasons preventing your company from moving faster with generative AI?



Sources: Bain Generative AI Surveys, October 2023 (N=198) and February 2024 (N=200)



What are the top reasons preventing your company from moving faster with generative AI?



Percentage-point change between October 2023 and February 2024

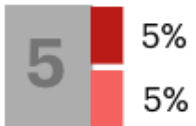
Sources: Bain Generative AI Surveys, October 2023 (N=198) and February 2024 (N=200)

Data Readiness for Generative AI



Tech Nontech

Data readiness for generative AI



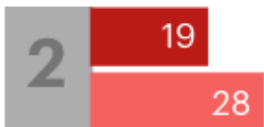
Structured and unstructured data are centralized, quality controlled and readily available



Structured and unstructured data are available, but both need to be quality controlled and centralized



Structured data is readily available, unstructured data takes some time to collect, process and check quality



Structured and unstructured data can be made available, but doing so requires significant effort and enhancements



Structured and unstructured data are not readily available.

Source: Bain Generative Artificial Intelligence Survey, February 2024 (N=200)

Data Readiness for Generative AI



Tech Nontech

Data readiness for generative AI



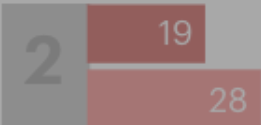
Structured and unstructured data are centralized, quality controlled and readily available



Structured and unstructured data are available, but both need to be quality controlled and centralized



Structured data is readily available, unstructured data takes some time to collect, process and check quality



Structured and unstructured data can be made available, but doing so requires significant effort and enhancements



Structured and unstructured data are not readily available.

Source: Bain Generative Artificial Intelligence Survey, February 2024 (N=200)

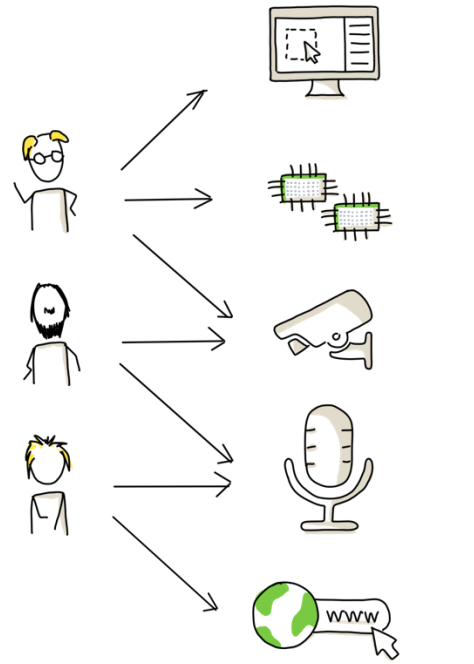


Data Management

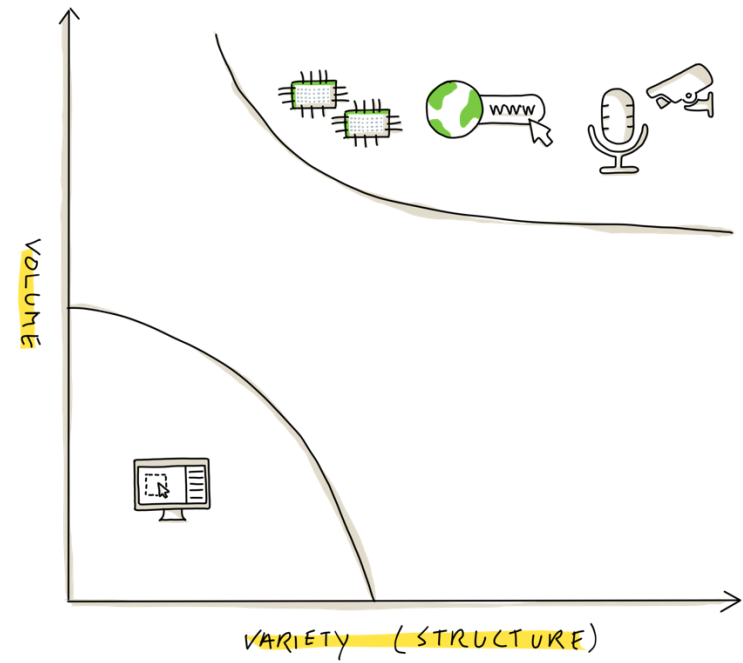
- Defining Data
- Data Producers
- Big Data Vs
- The Bigger Picture
- **Data Management**
- Data Technology



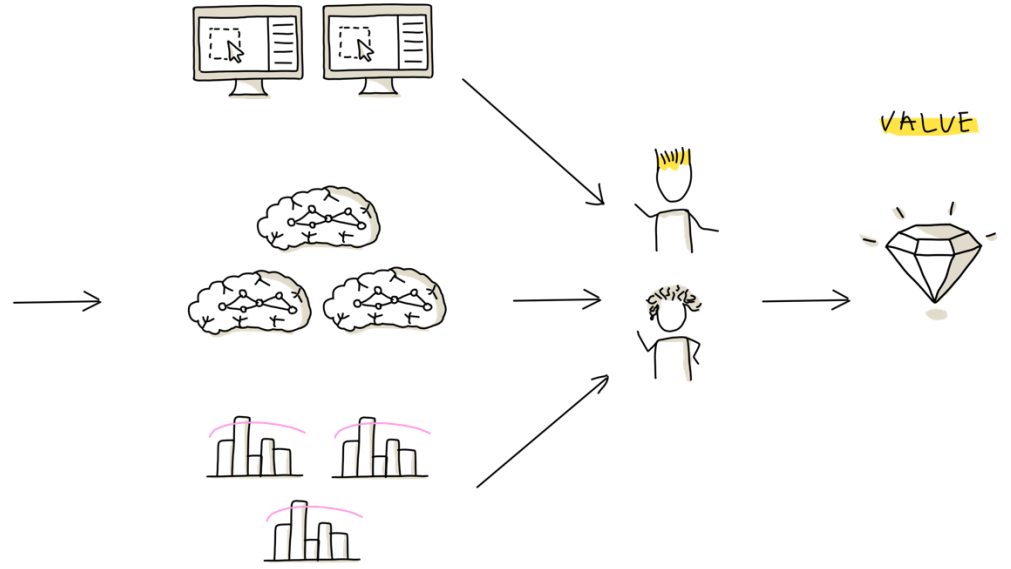
DATA PRODUCERS (SOURCES)



DATA

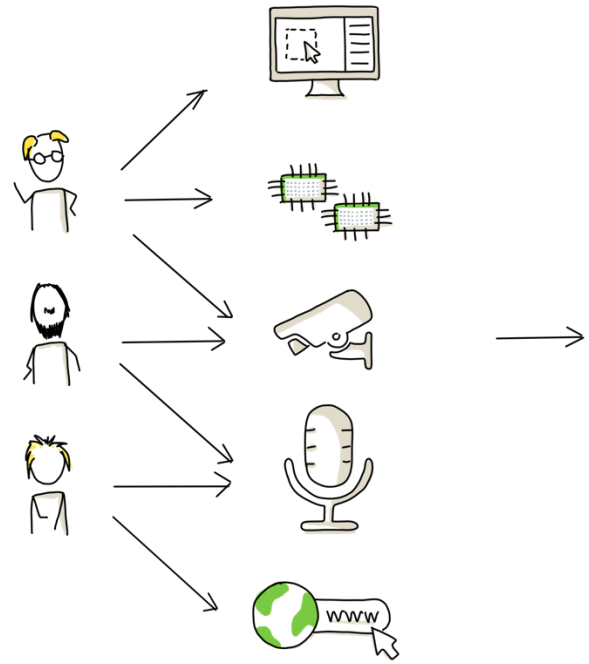


DATA CONSUMERS

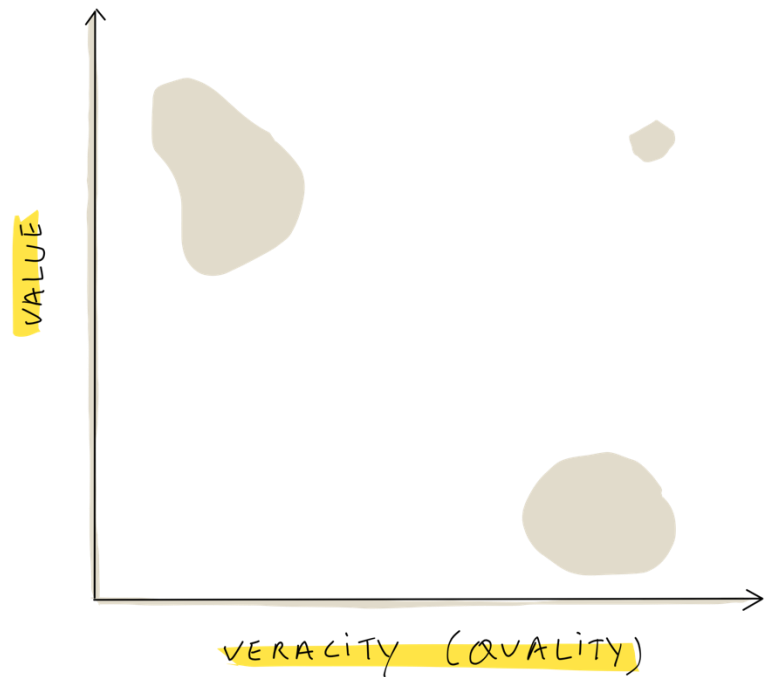




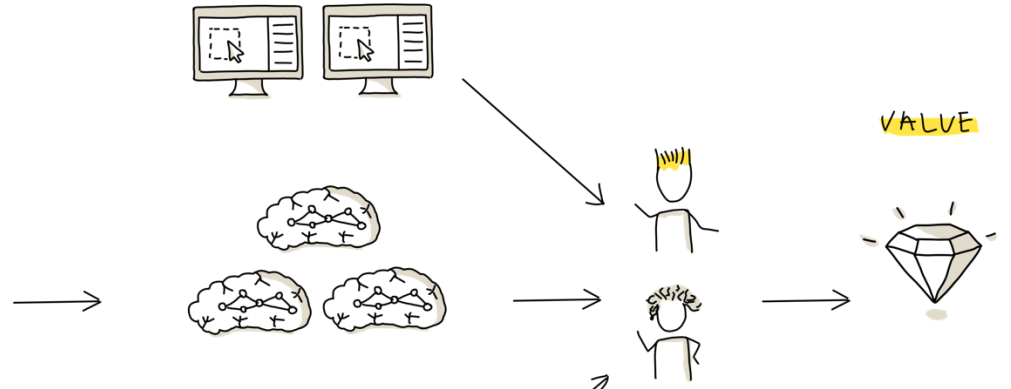
DATA PRODUCERS (SOURCES)



DATA

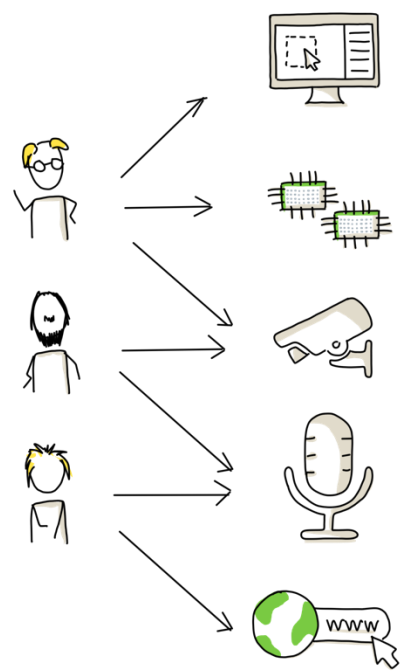


DATA CONSUMERS

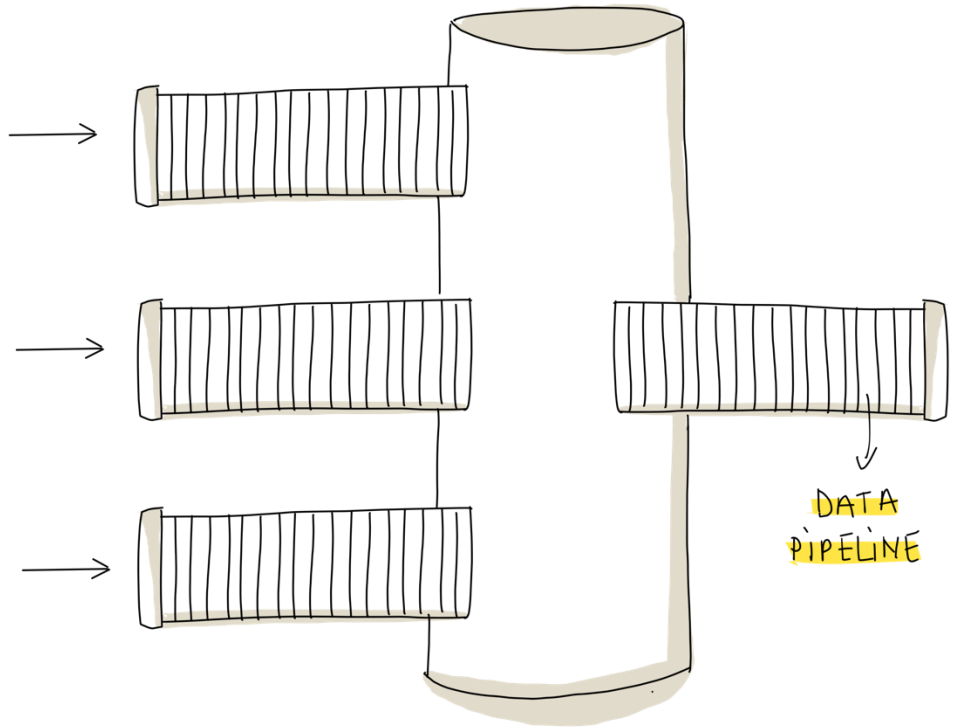




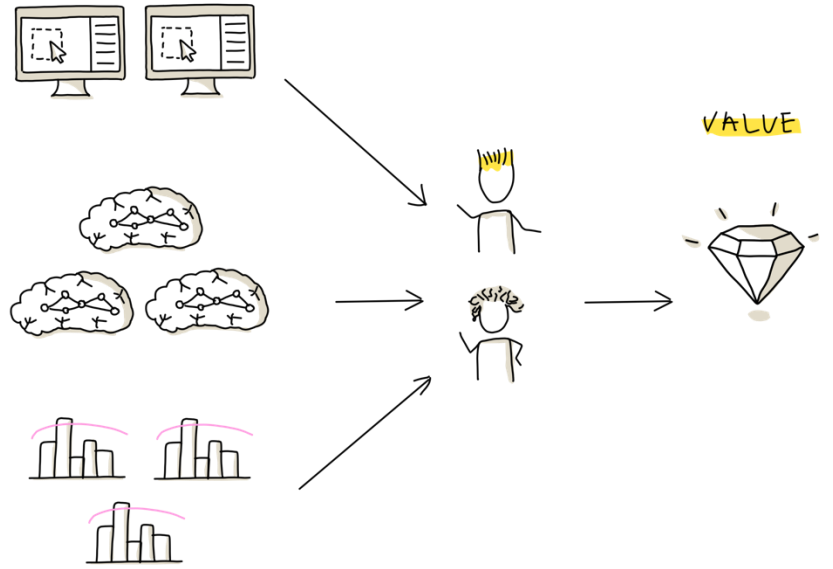
DATA PRODUCERS
(SOURCES)



DATA PLATFORM

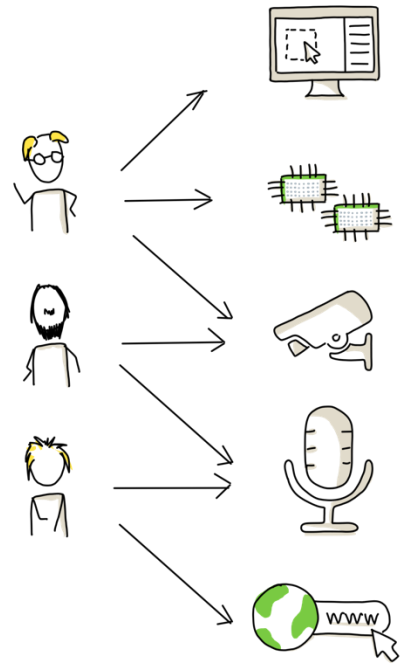


DATA CONSUMERS

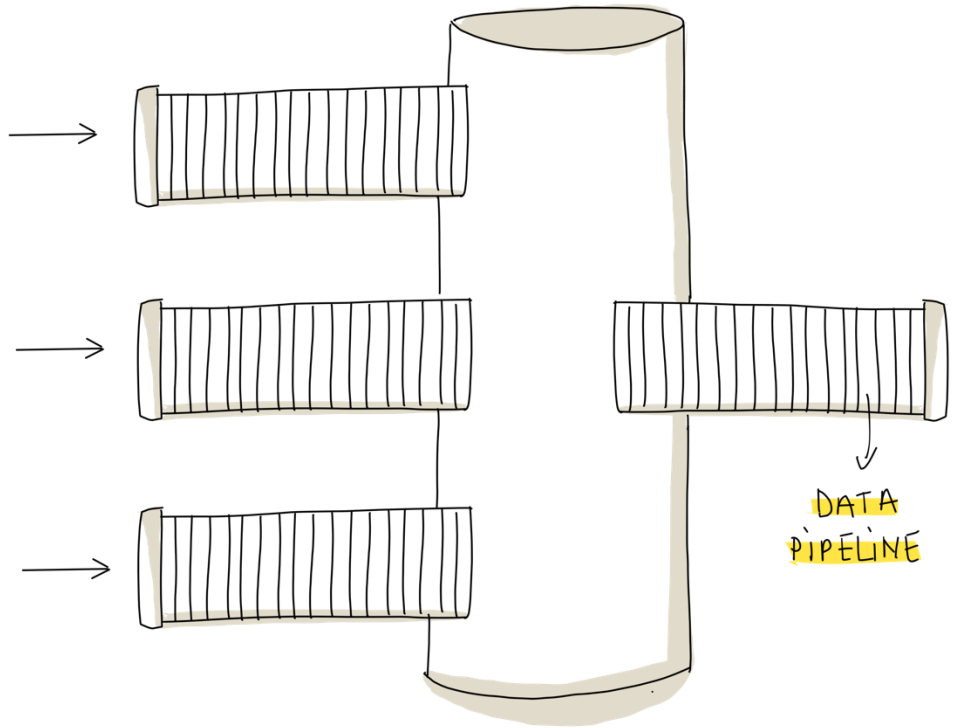




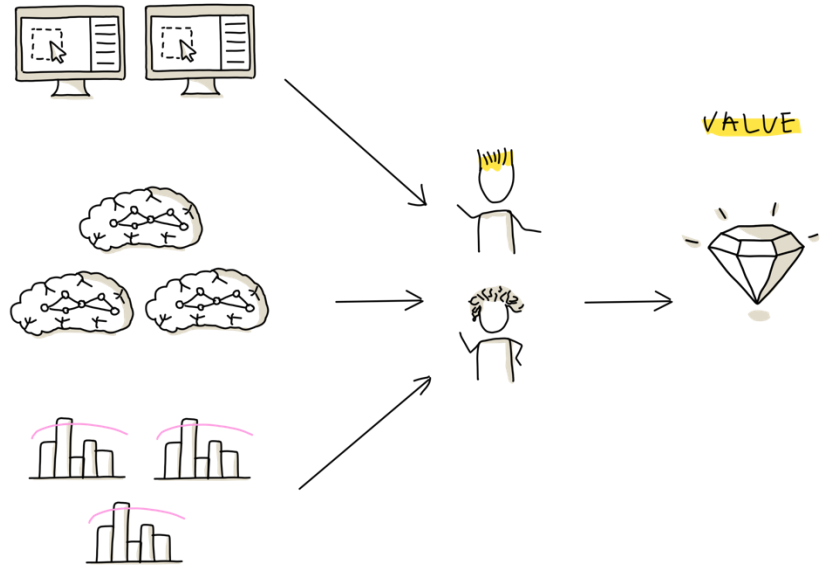
DATA PRODUCERS
(SOURCES)



DATA PLATFORM

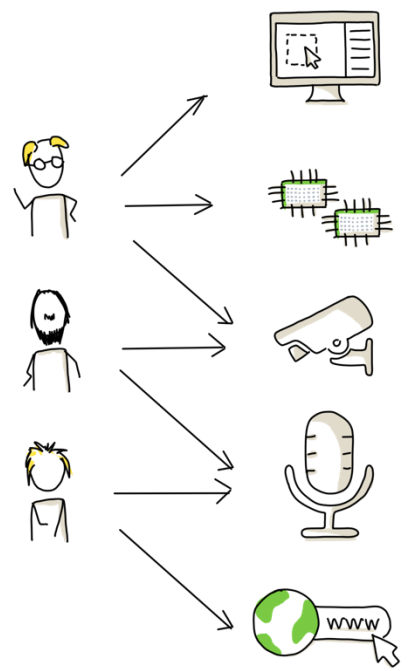


DATA CONSUMERS

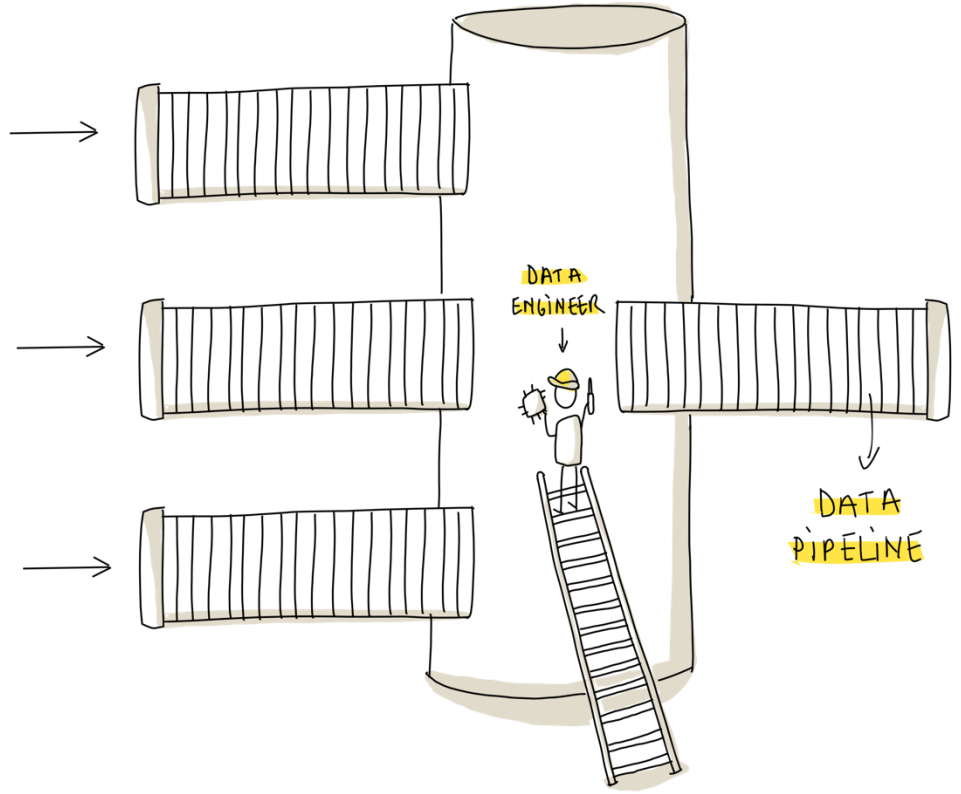




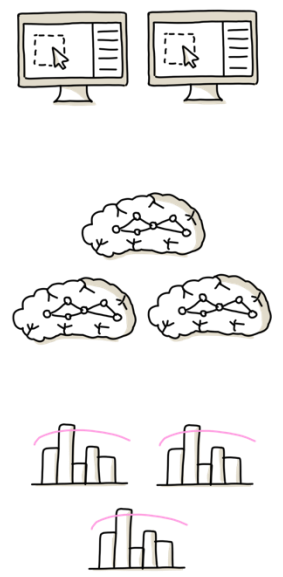
DATA PRODUCERS (SOURCES)



DATA PLATFORM



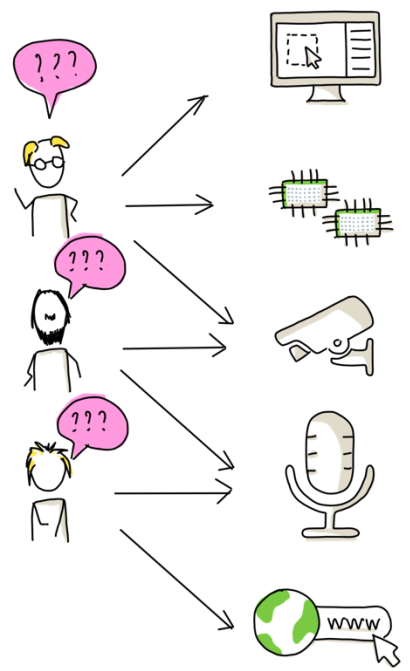
DATA CONSUMERS



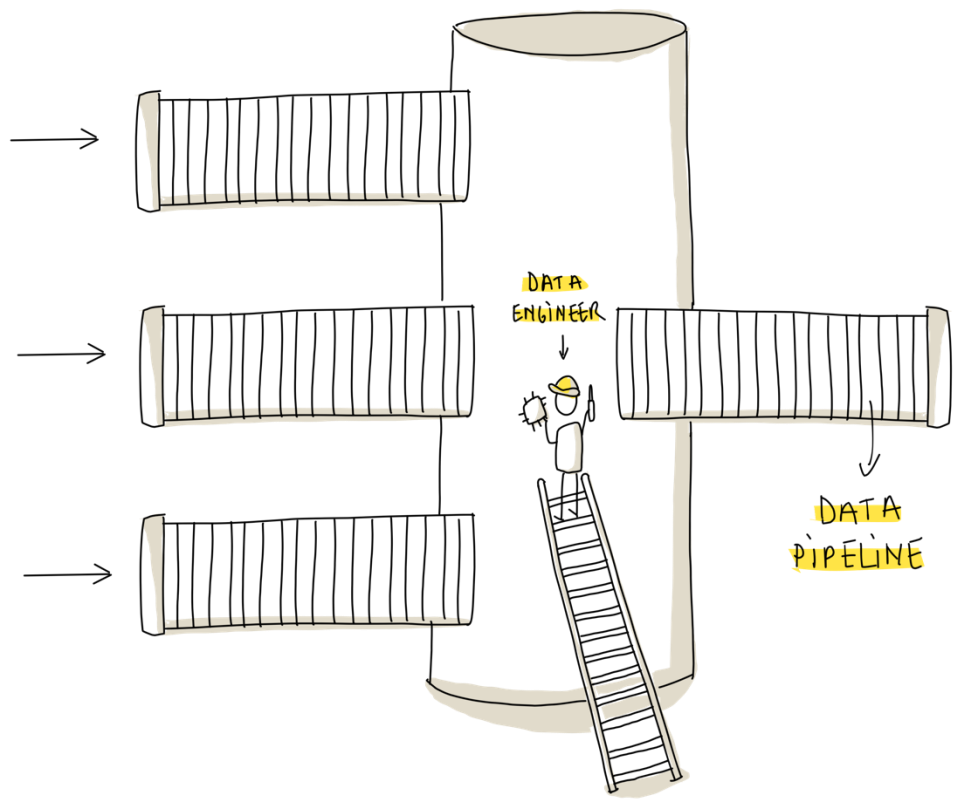
DATA ENGINEERING



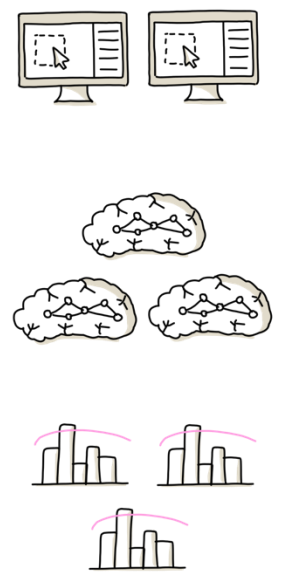
DATA PRODUCERS (SOURCES)



DATA PLATFORM



DATA CONSUMERS



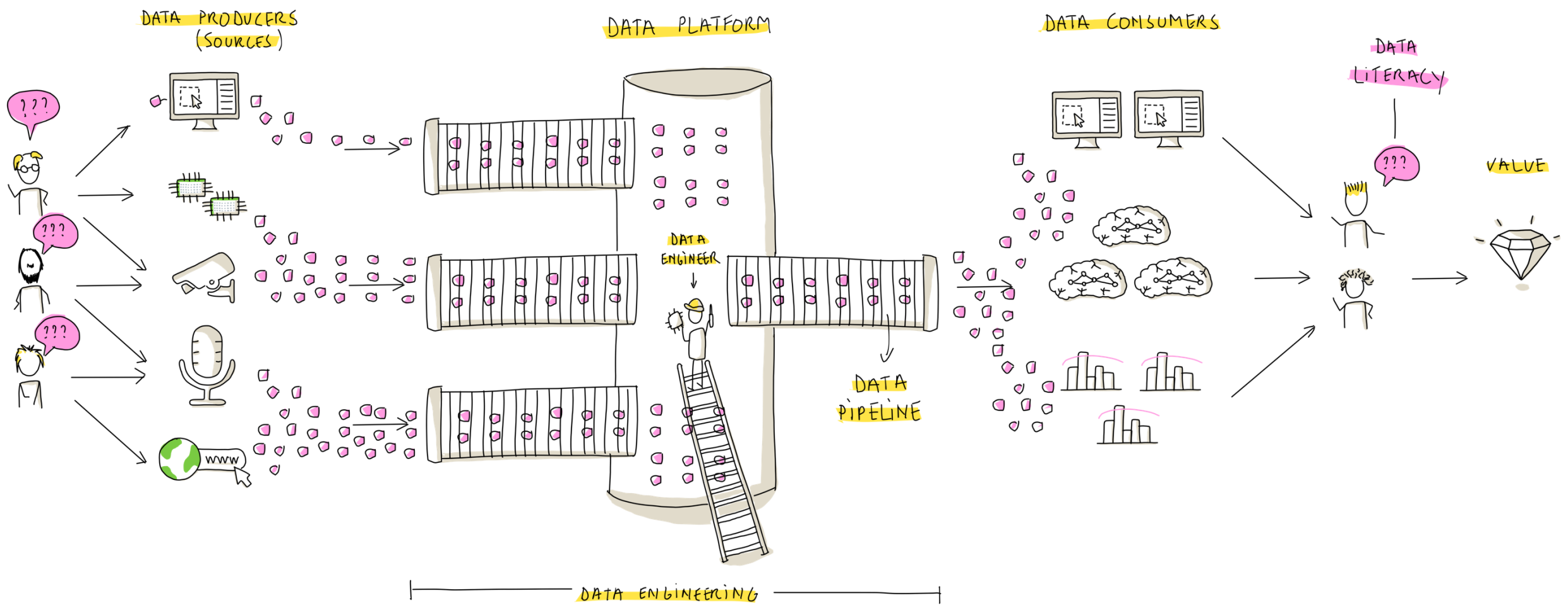
DATA LITERACY



VALUE



DATA ENGINEERING





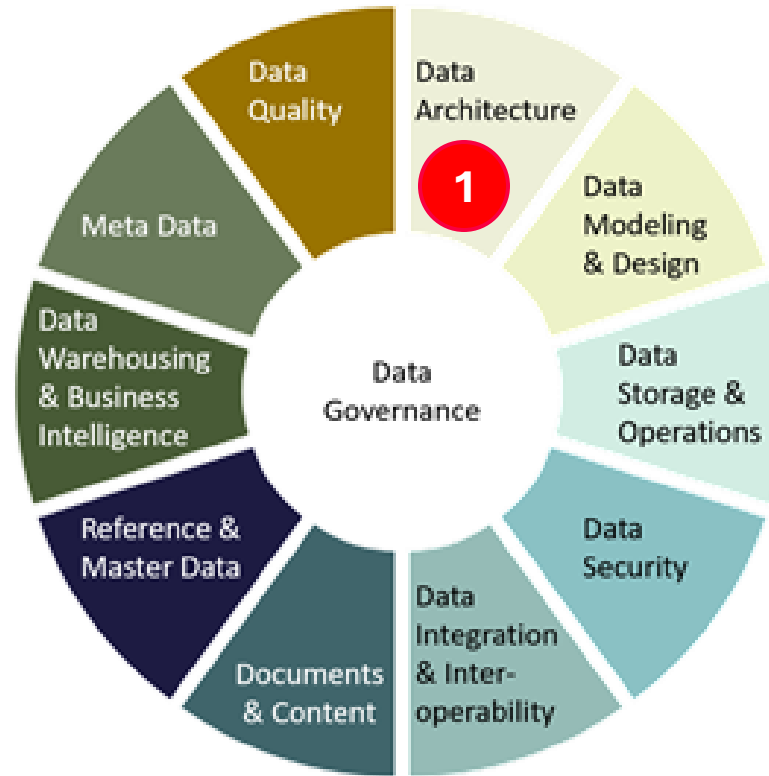
There is **Room for**
Improvement...

DMBOK – DAta MAnagement WHEEL



Data Governance defines the rules of play, Data Management executes these rules

Data Architecture



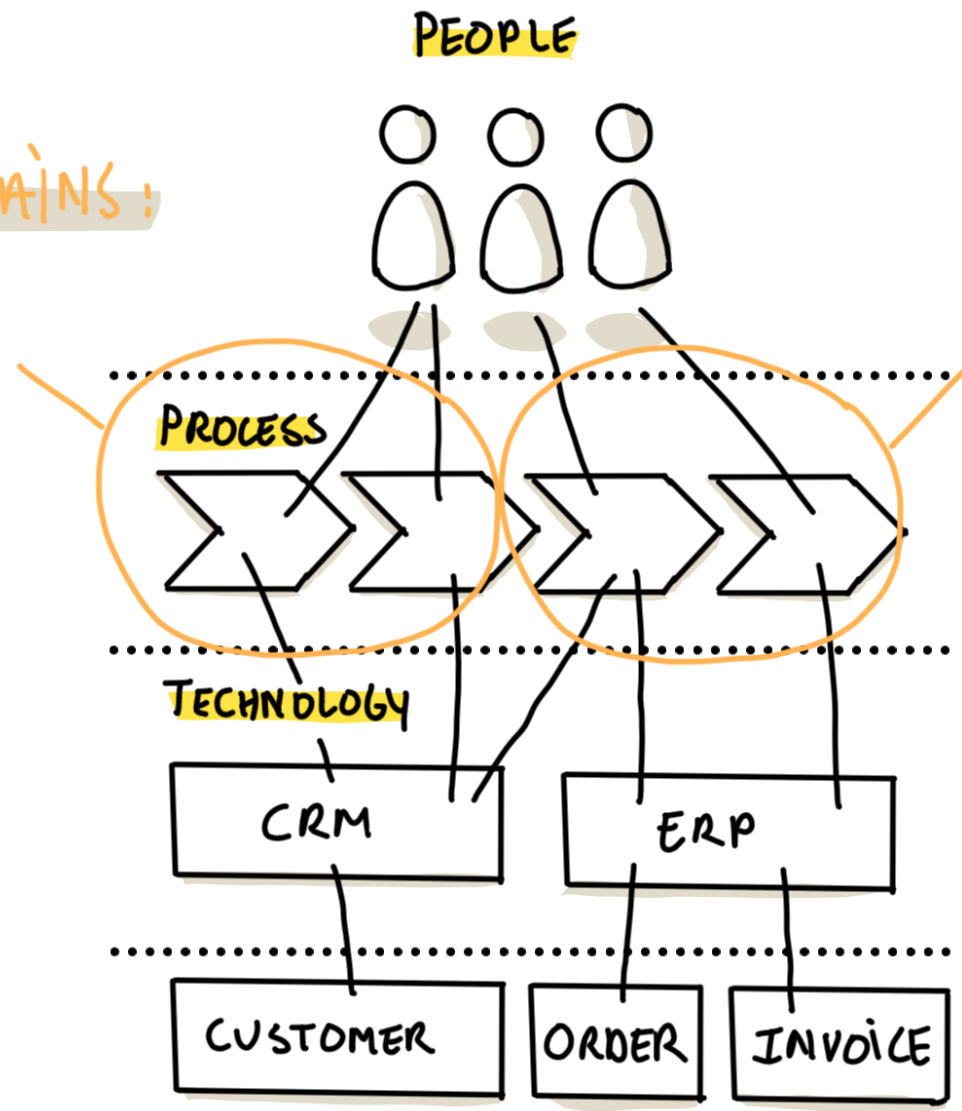
Which data do we have and how does it flow through the organization?



BUS. DOMAINS:

SALES

FINANCE/
ACCOUNTING



DATA (ENTITIES)



Business Domain

A specific **area of expertise or knowledge within a company** that focuses on particular business functions or industries. For example, "Finance" or "Human Resources" are business domains. It defines the scope of responsibilities, rules, and activities in a specific field.

Business Process

A series of tasks or steps carried out to achieve a specific business goal. **It's how work gets done within an organization.** For example, the **order-to-cash process** includes all the steps from receiving an order to collecting payment.



Entities

An entity is a **conceptual representation of something relevant to a business**, usually a tangible or abstract entity, **that groups data**. Examples:

- **Customer**, grouping data like name, address, and purchase history
- **Product**, grouping data like product ID, price, and stock level.

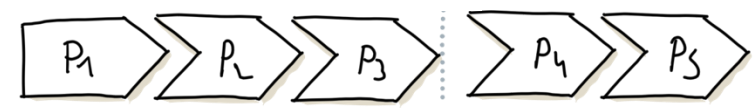


Domains, Processes and Entities

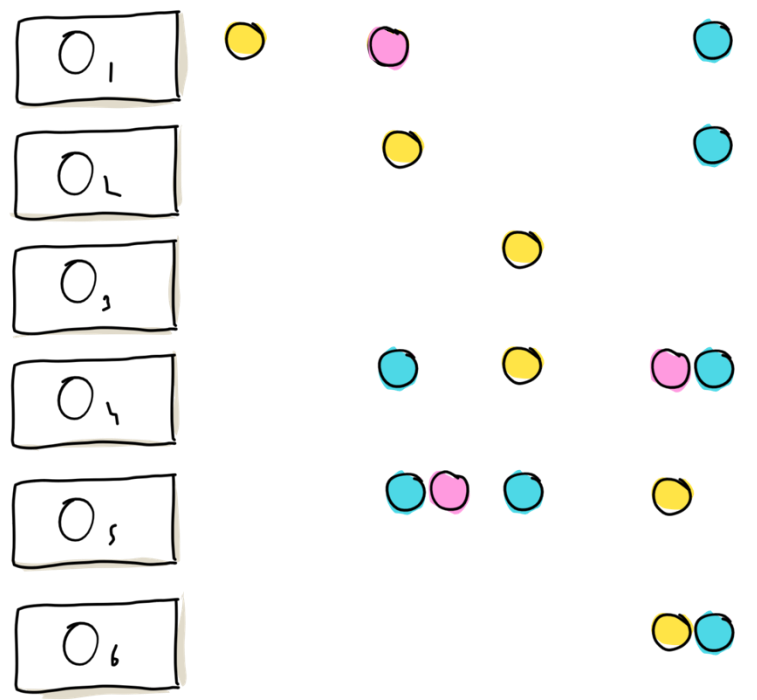
BUSINESS DOMAINS



BUSINESS PROCESSES

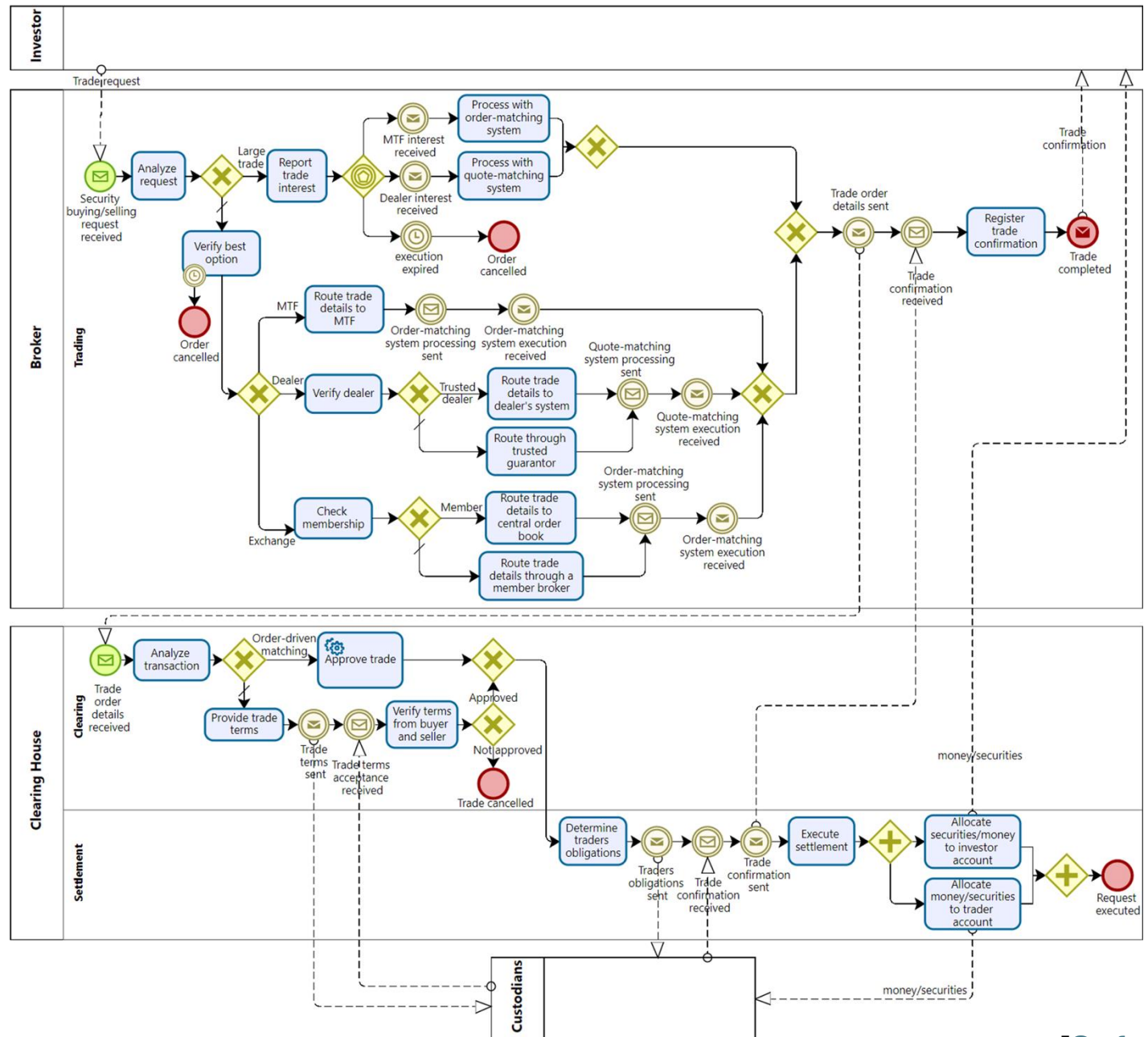


DATA ENTITIES



● = CREATE ● = USE ● = UPDATE

“Process **P1** (part of Domain **DA**) creates Entity **O1**. This entity is updated in Process **P2** and used in Process **P4**.”

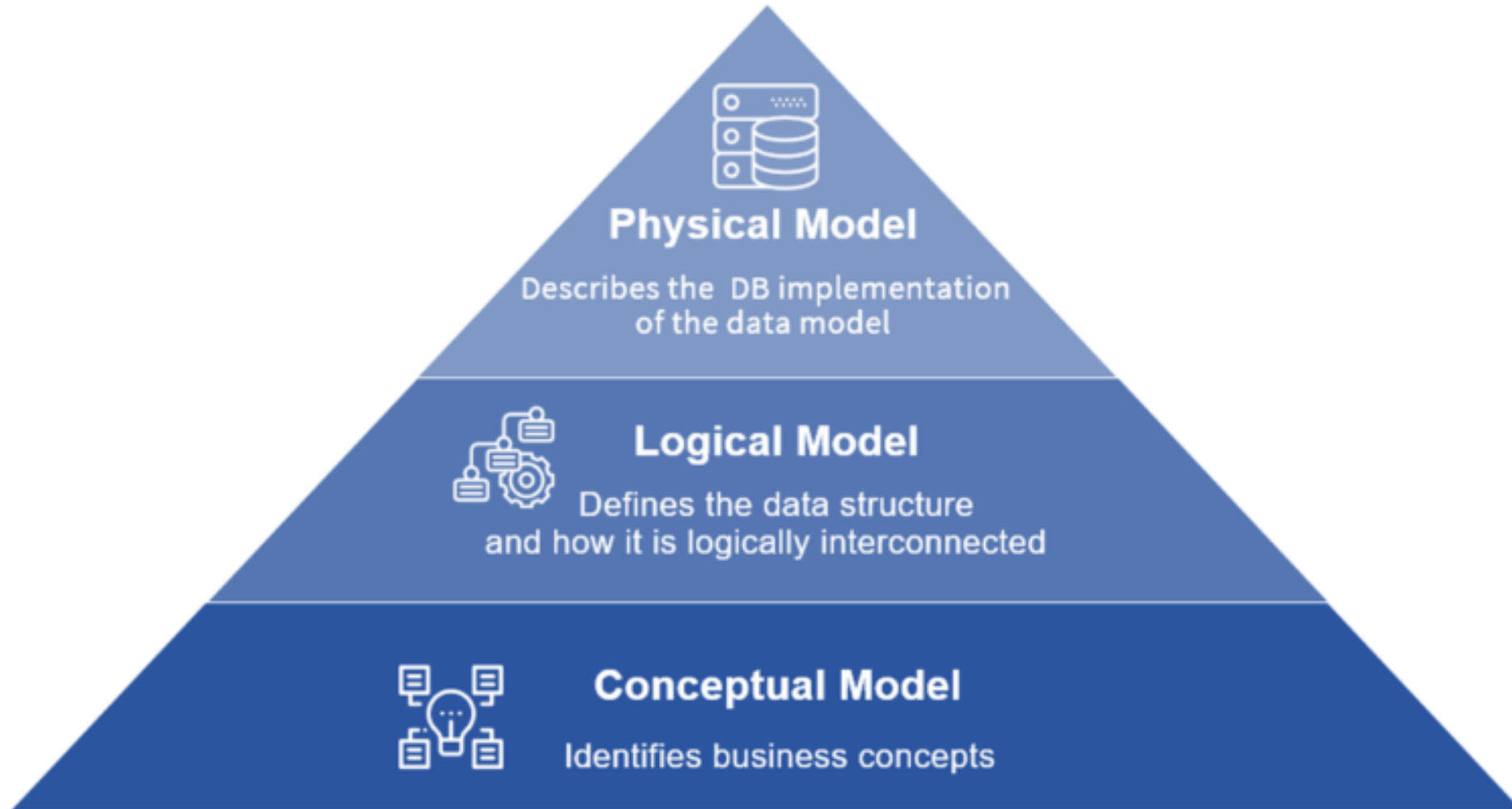


Data Modeling & Design



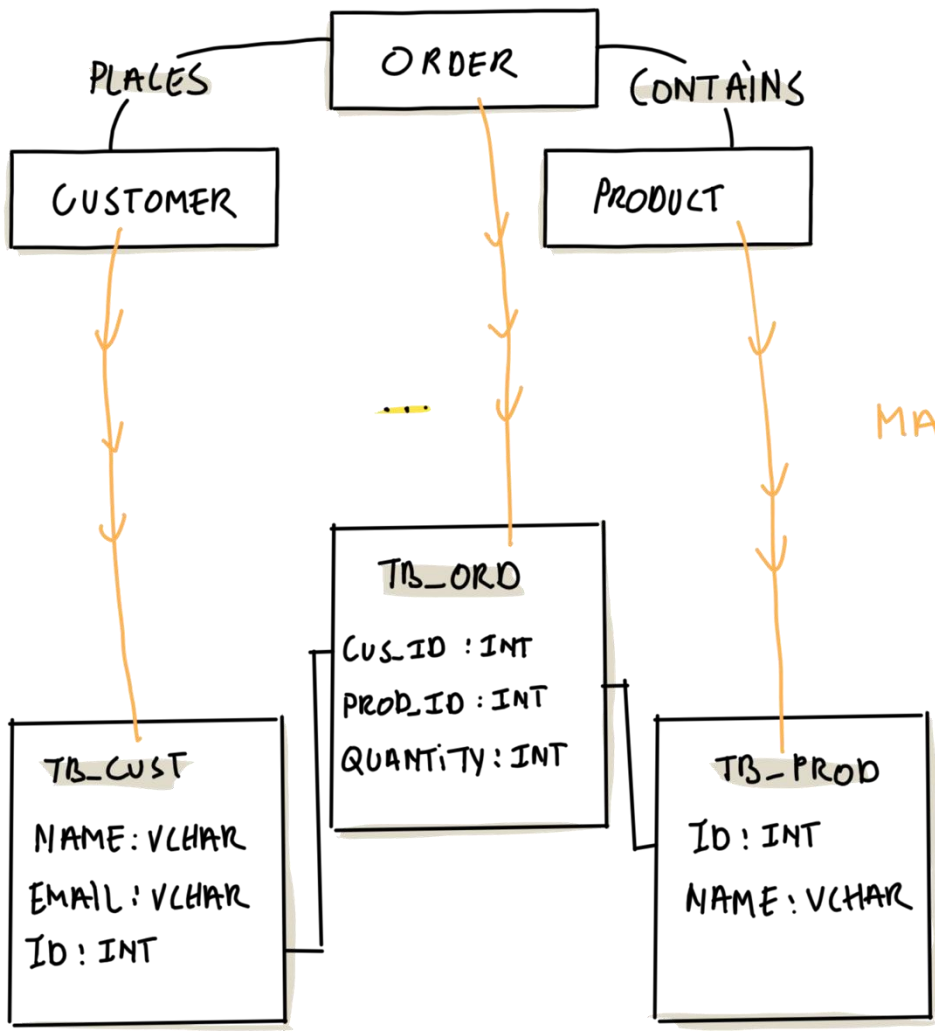
Providing a Common Vocabulary
around Data – Often
communicated in a ‘Data Model’

Physical, Logical & Conceptual Model





BUSINESS VIEW



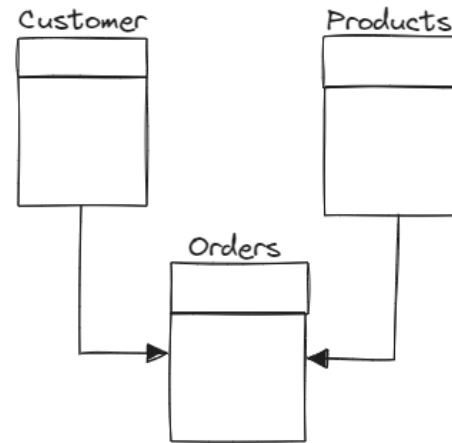
MAPPING

DATA (BASE) VIEW

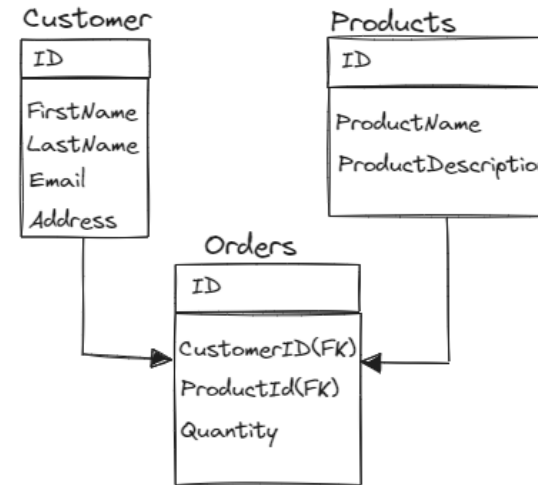
Physical, Logical & Conceptual Model



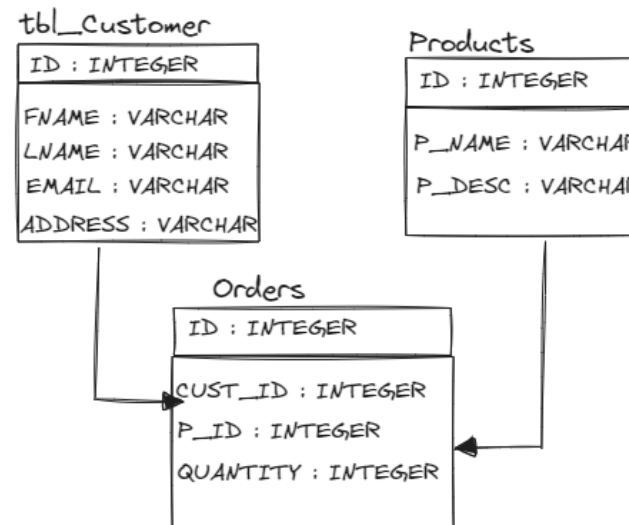
CONCEPTUAL MODEL

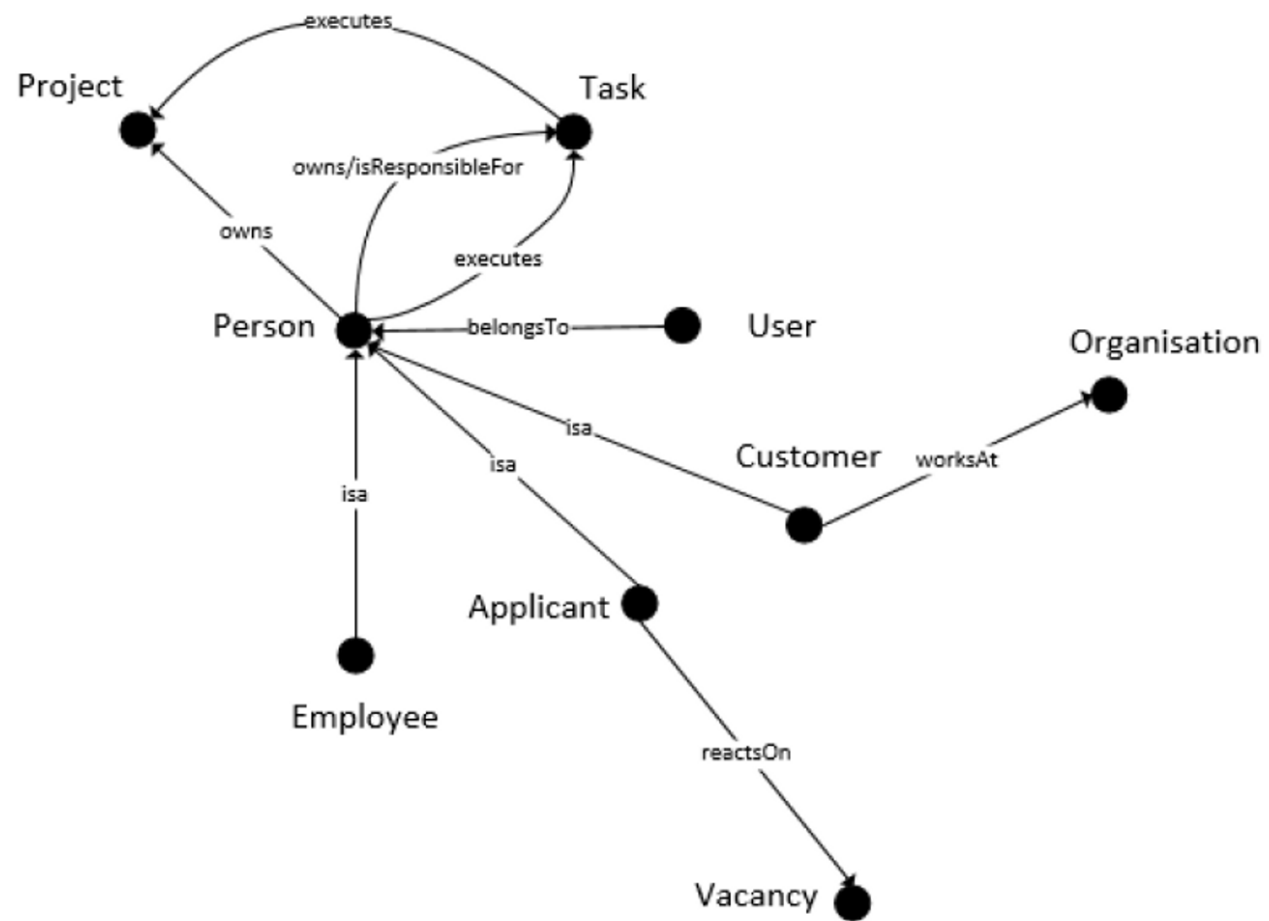


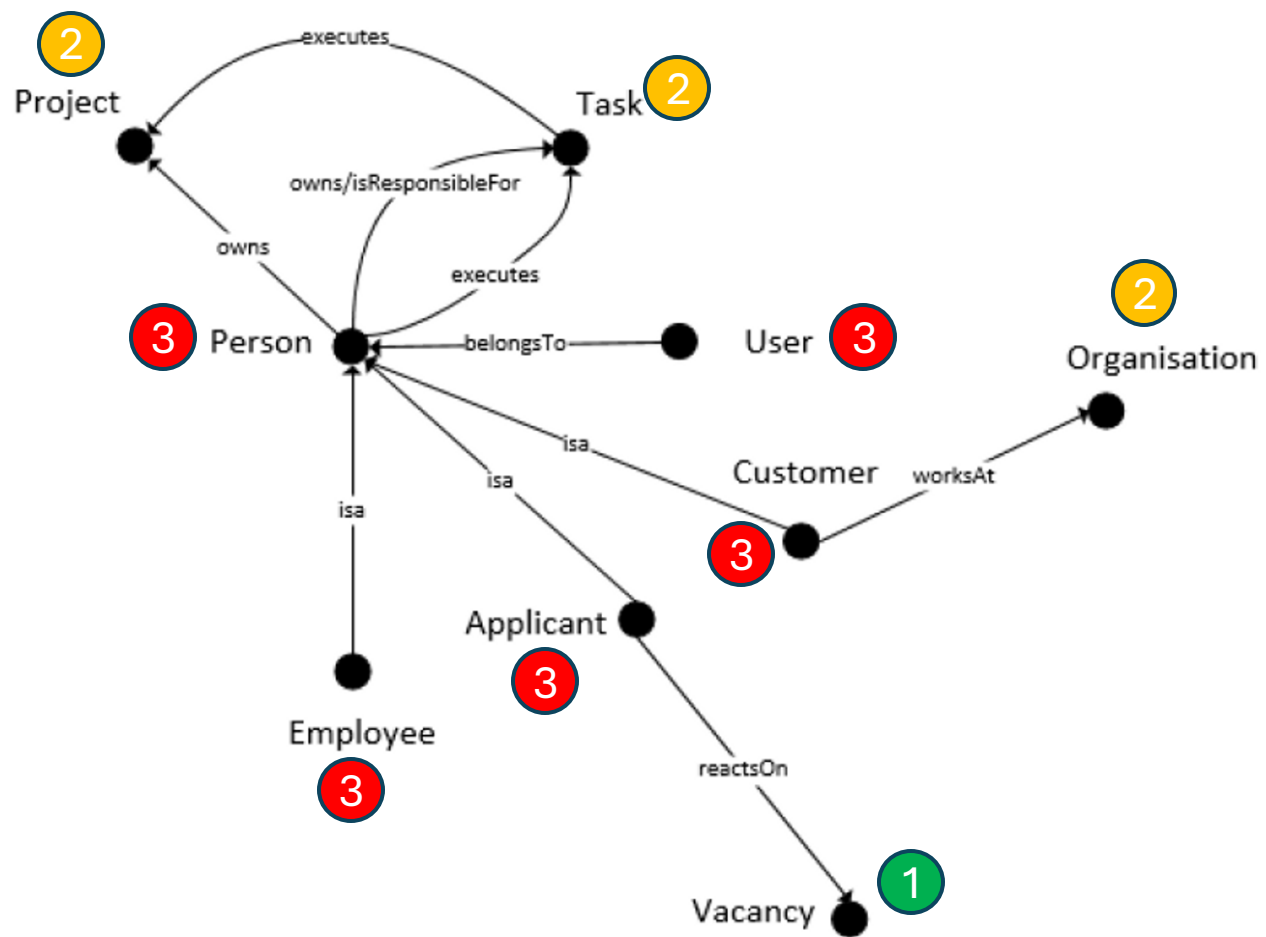
LOGICAL MODEL



PHYSICAL MODEL







1 Public data

2 Internal

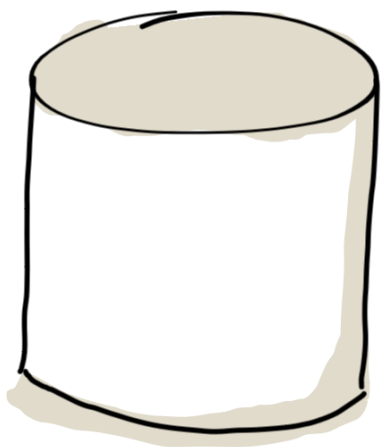
3 Sensitive

Classification: Adding Sensitivity Information!

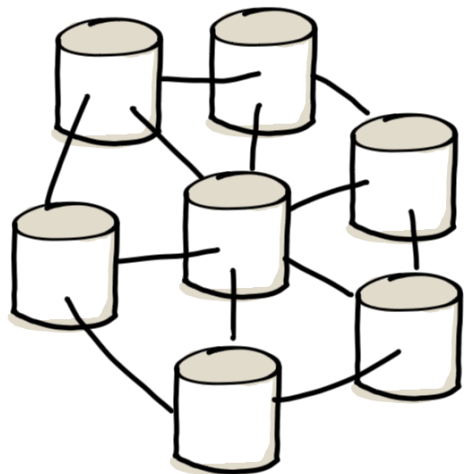
Data Storage & Operations



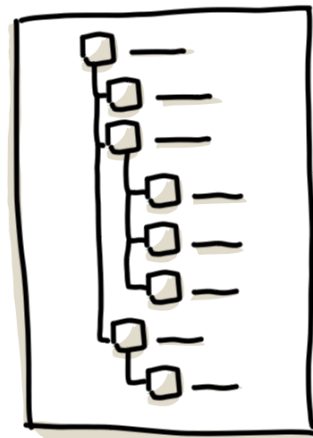
How do we store data for analytical use?



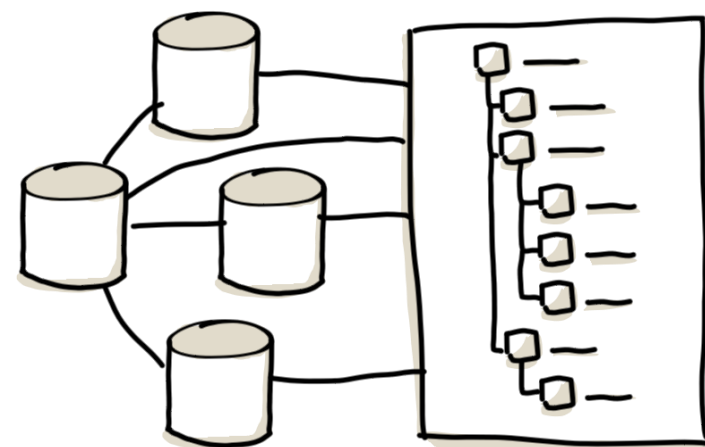
DATABASE



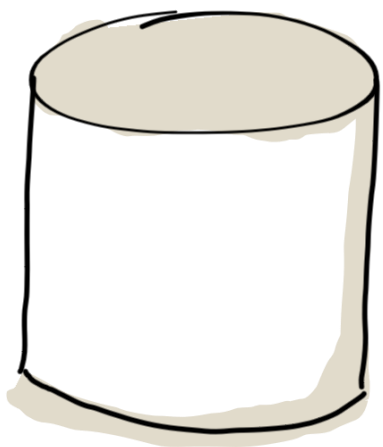
ANALYTICAL
DATABASE



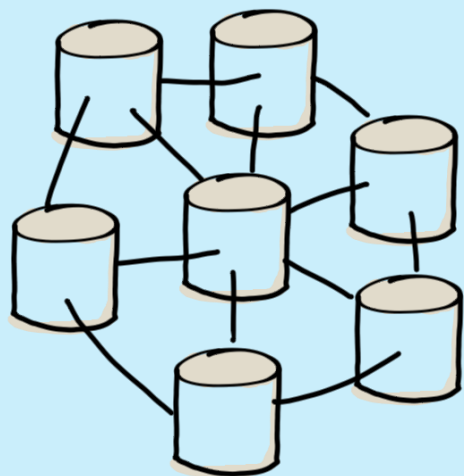
DATA LAKE



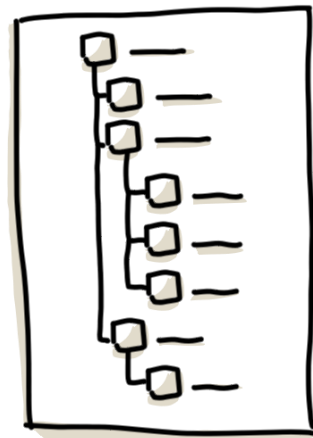
DATA
LAKEHOUSE



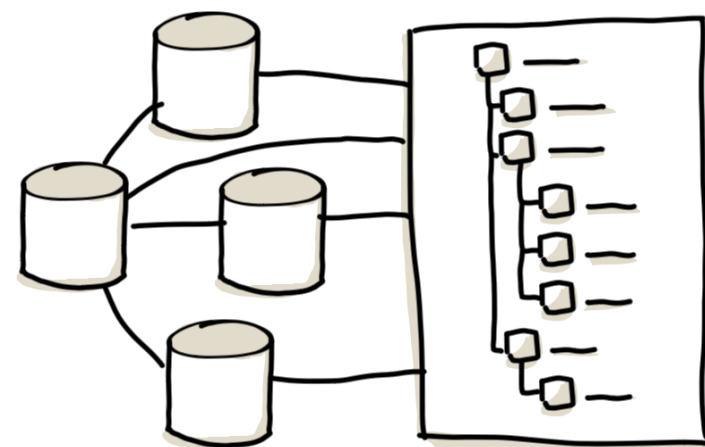
DATABASE



ANALYTICAL
DATABASE



DATA LAKE

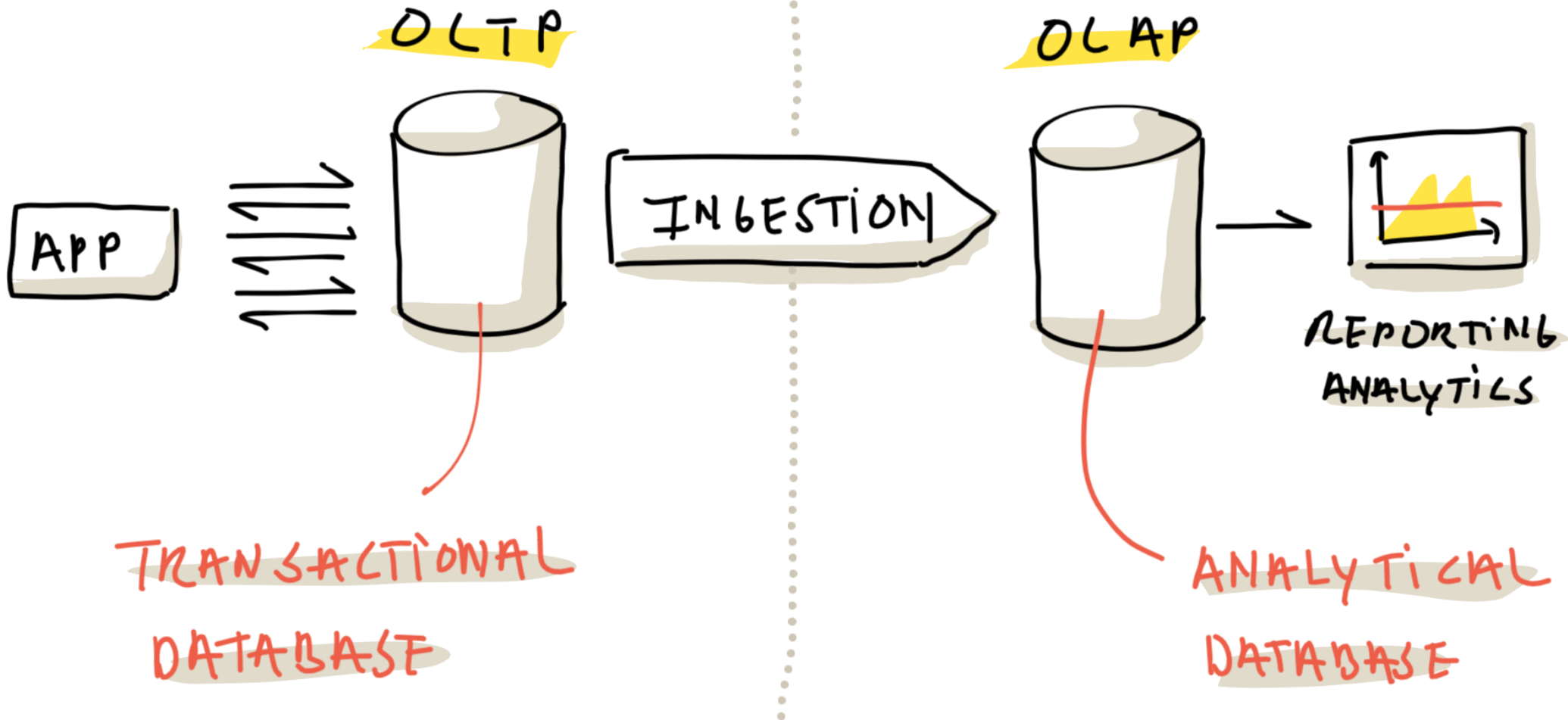


DATA
LAKEHOUSE



DATA SOURCES

DATA PLATFORM



OLTP

OLAP

INGESTION

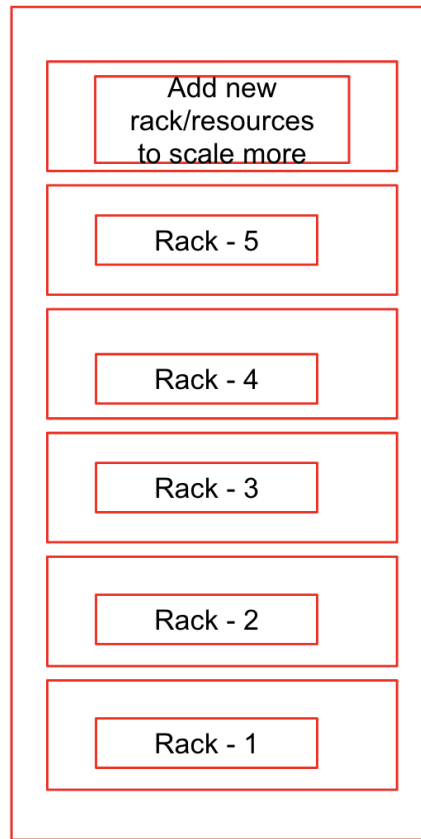
TRANSACTIONAL
DATABASE

ANALYTICAL
DATABASE

REPORTING
ANALYTICS



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

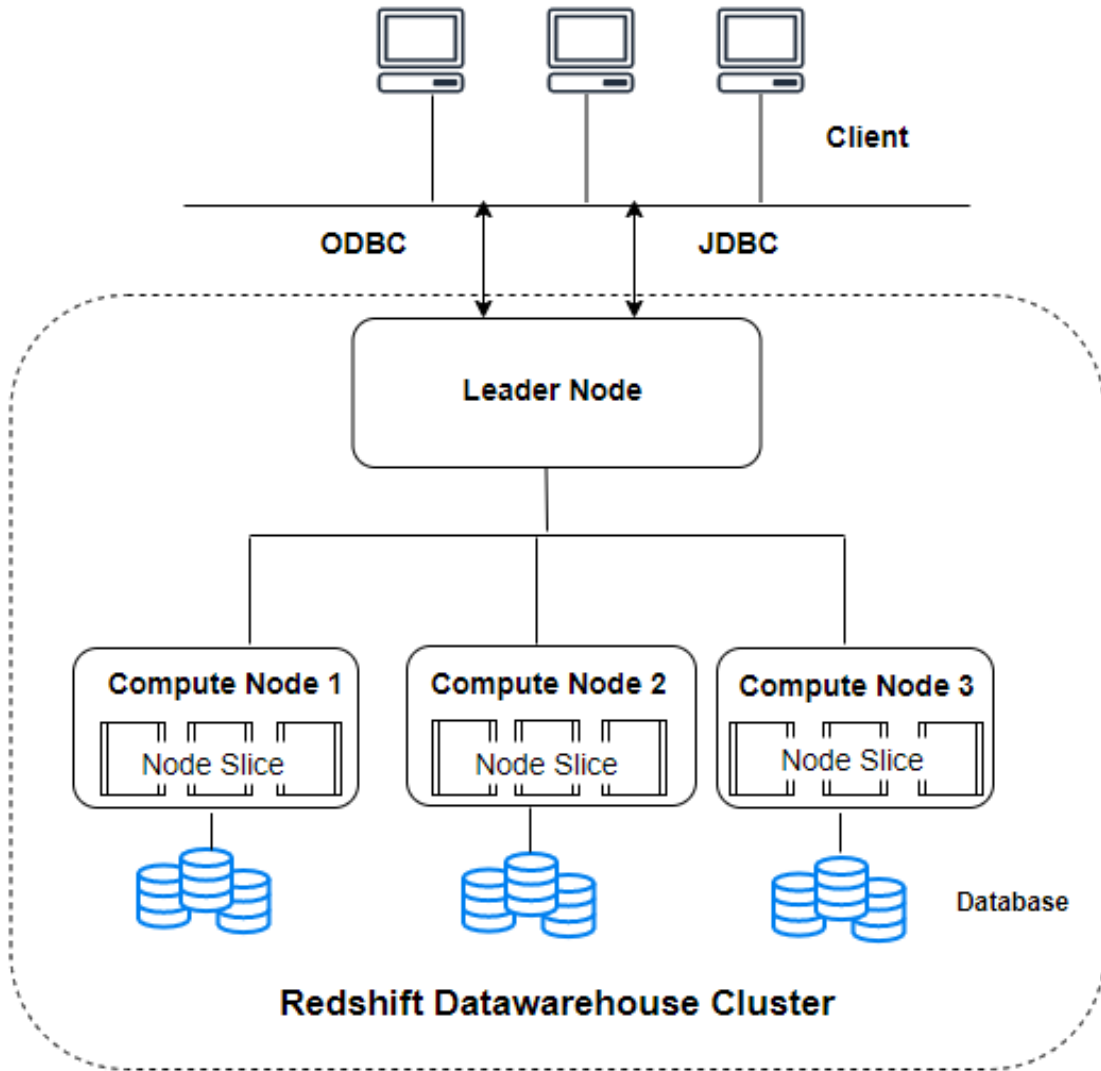
Host 2
192.168.1.2

Host 3
192.168.1.3

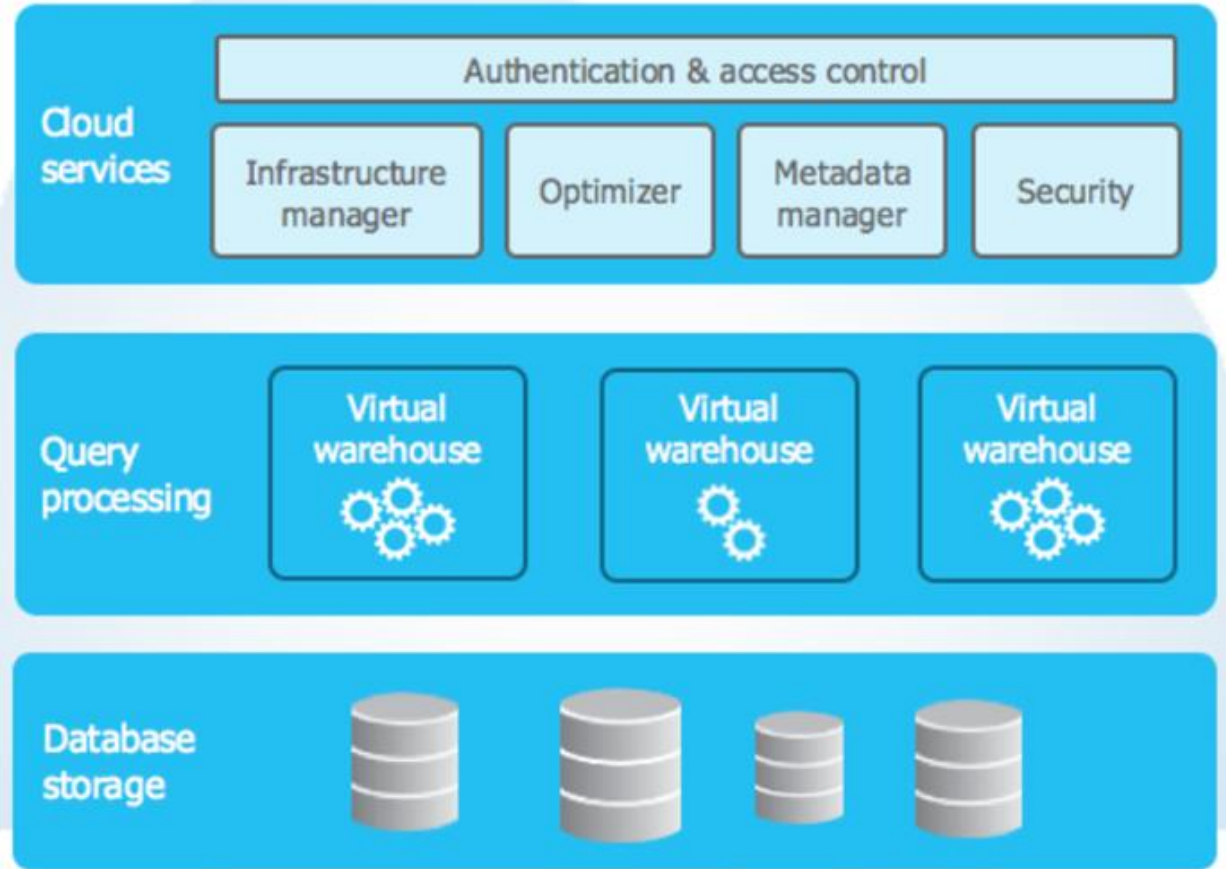
Host x
192.168.1.x

Add x+1 host to scale out

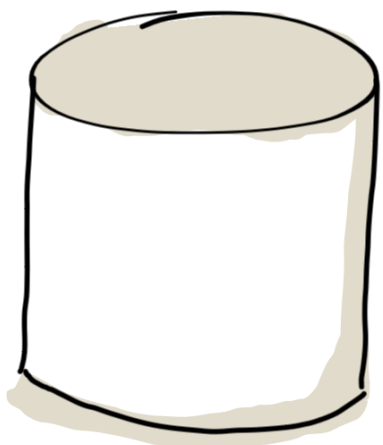
Example: Amazon Redshift



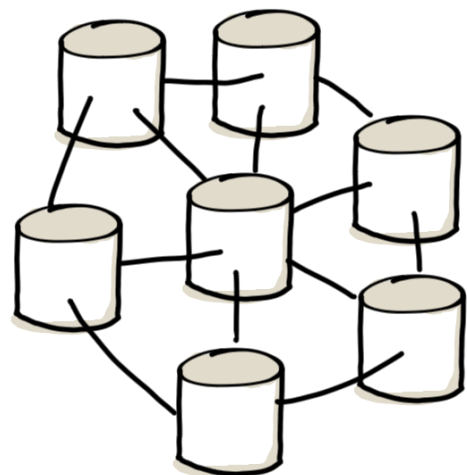
Example: Snowflake



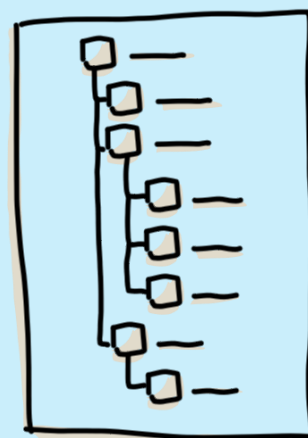




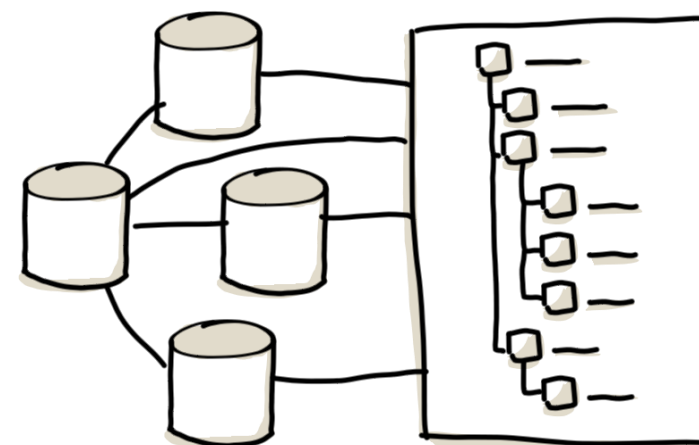
DATABASE



ANALYTICAL
DATABASE



DATA LAKE



DATA
LAKEHOUSE

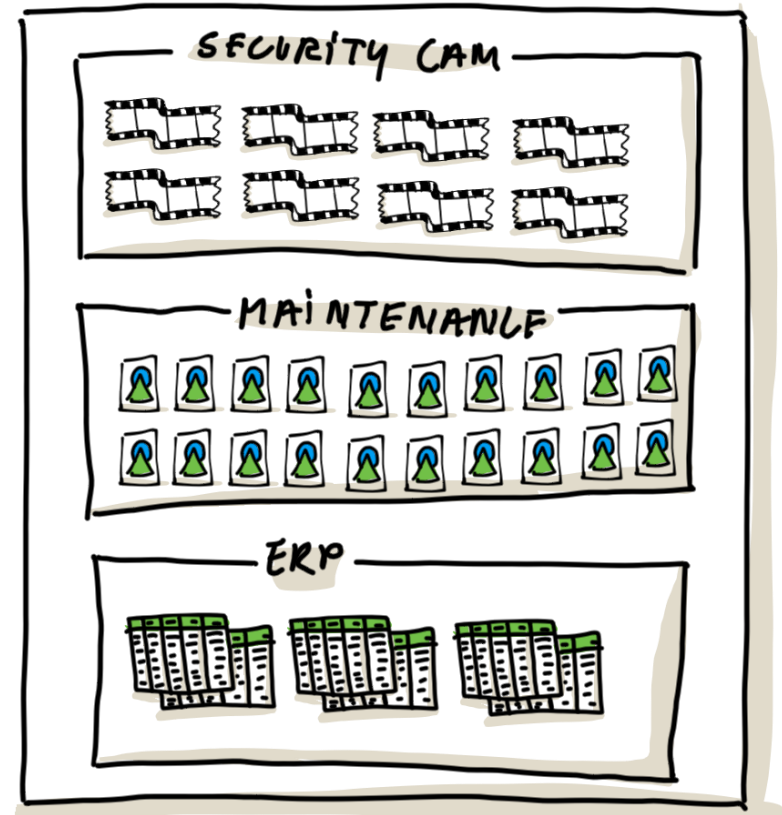


Data Lake

DATA SOURCES



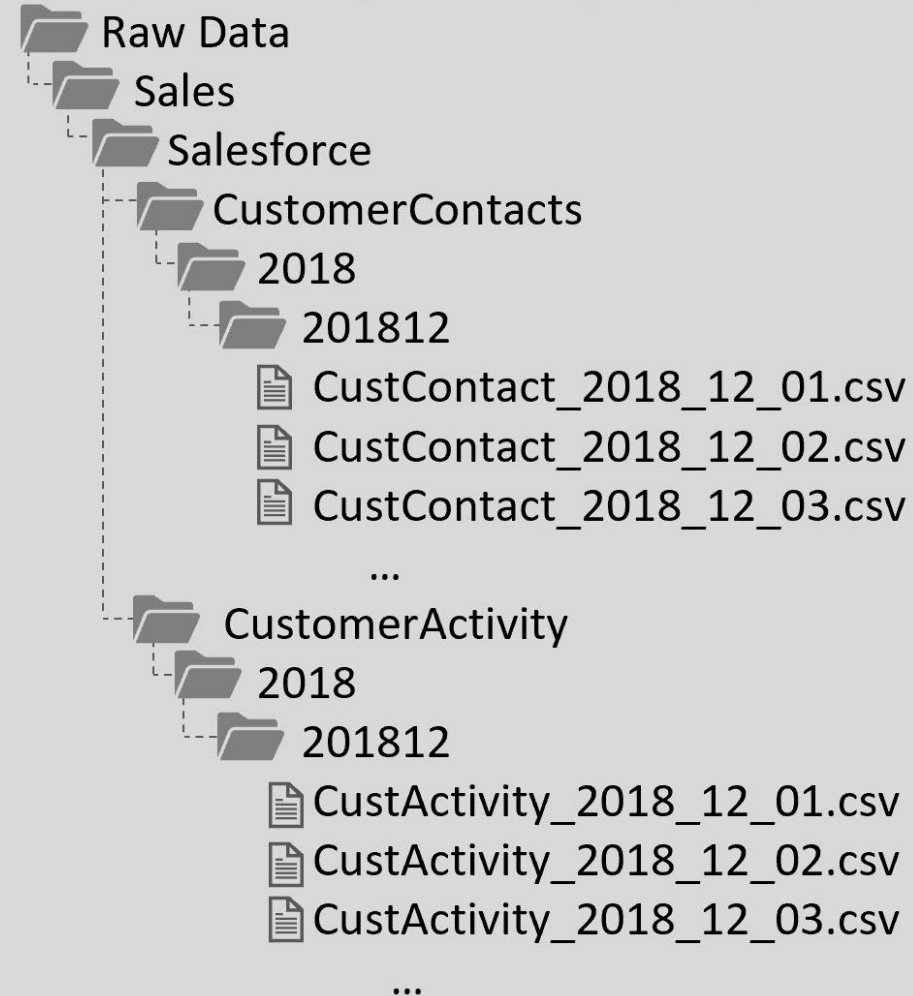
DATA LAKE

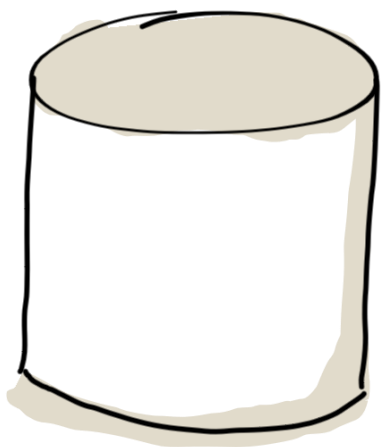


DATA PRODUCTS

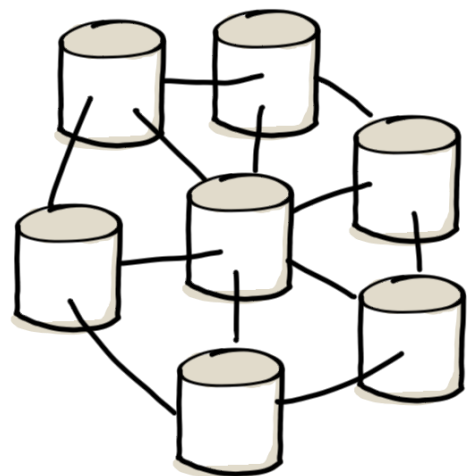


Data Lake: A **Structured** File Repository

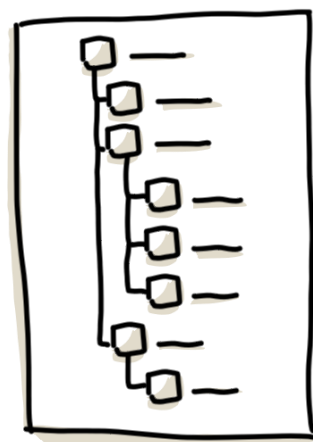




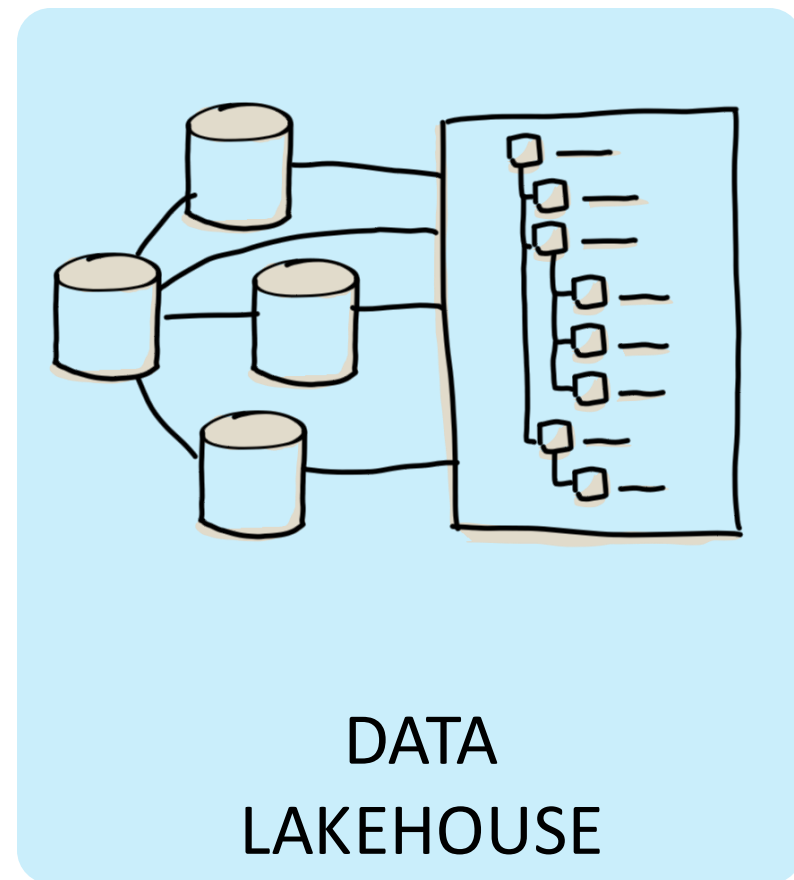
DATABASE



ANALYTICAL
DATABASE



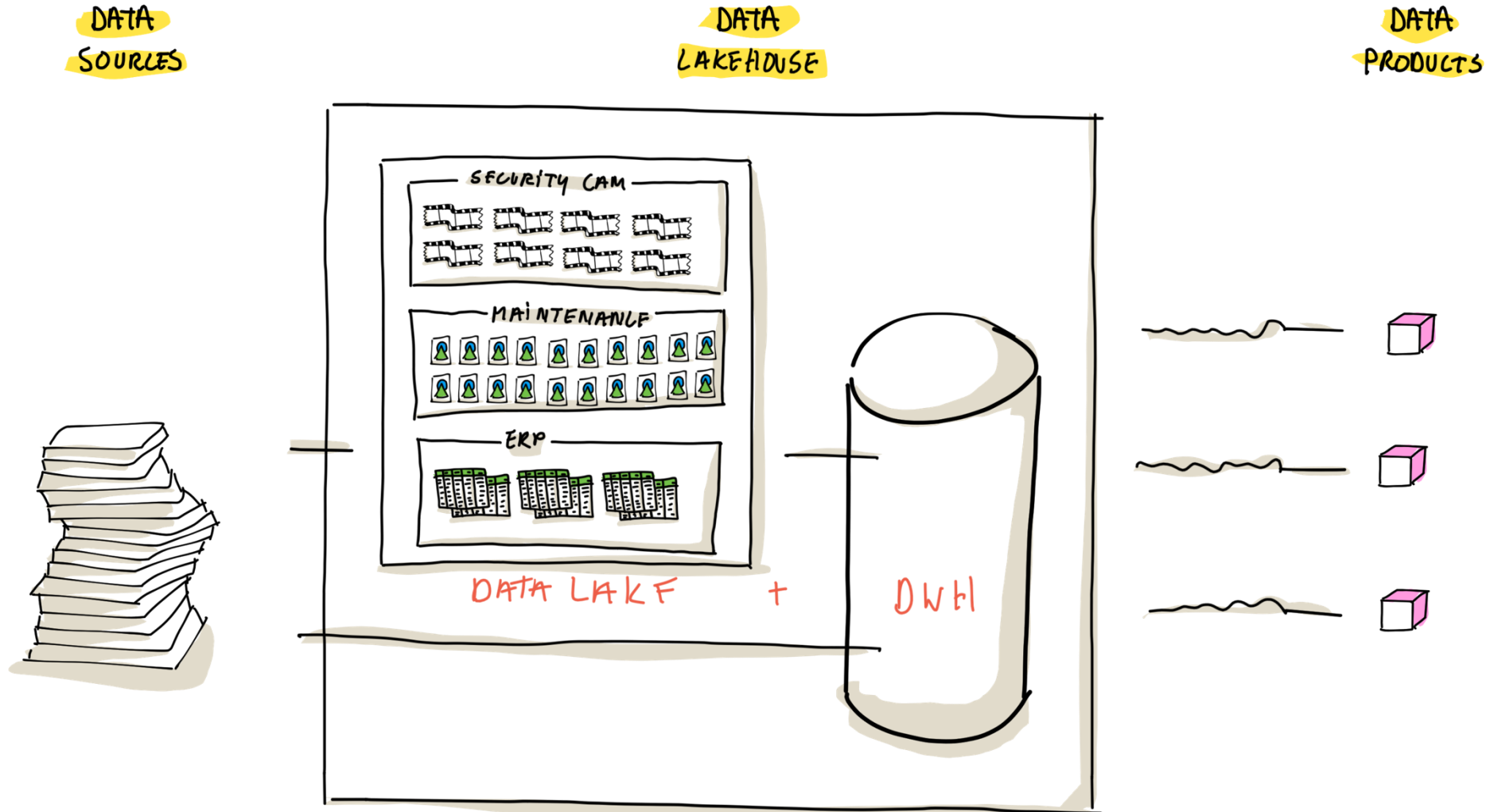
DATA LAKE



DATA
LAKEHOUSE



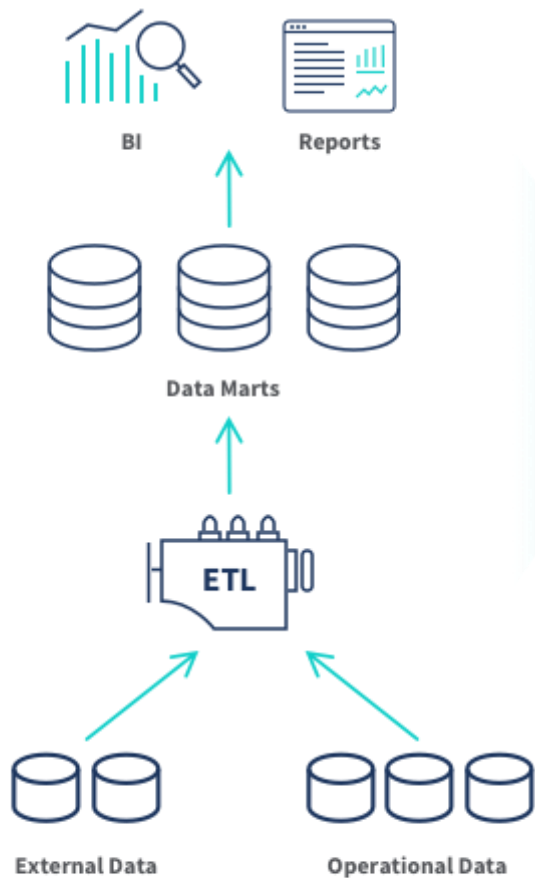
Data Lakehouse





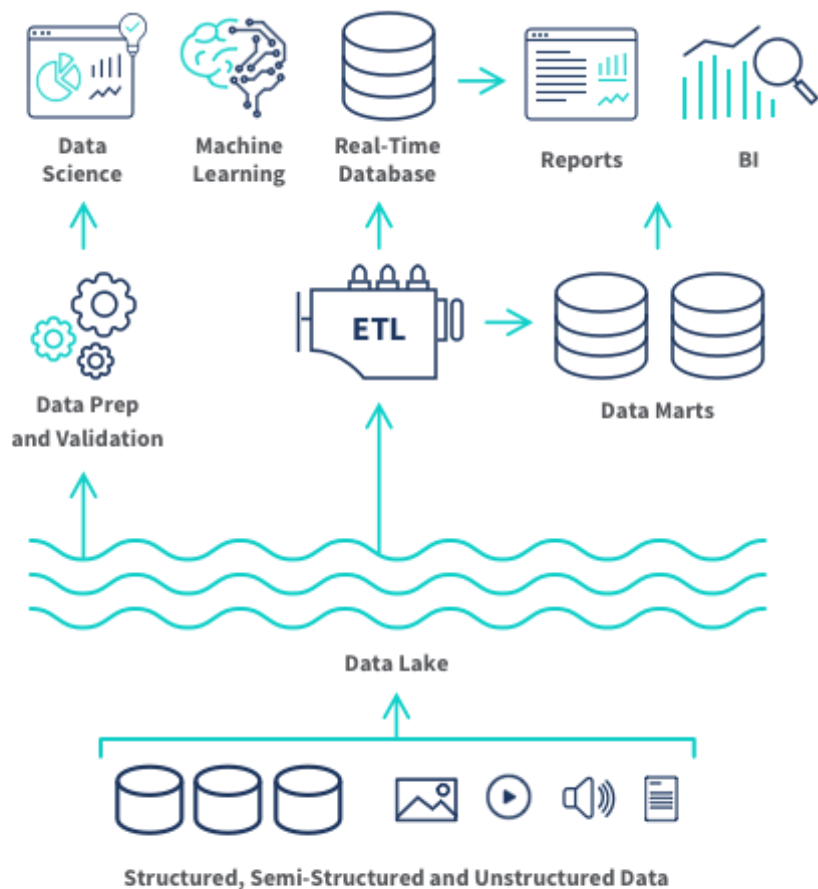
LATE 1980'S

Data Warehouse



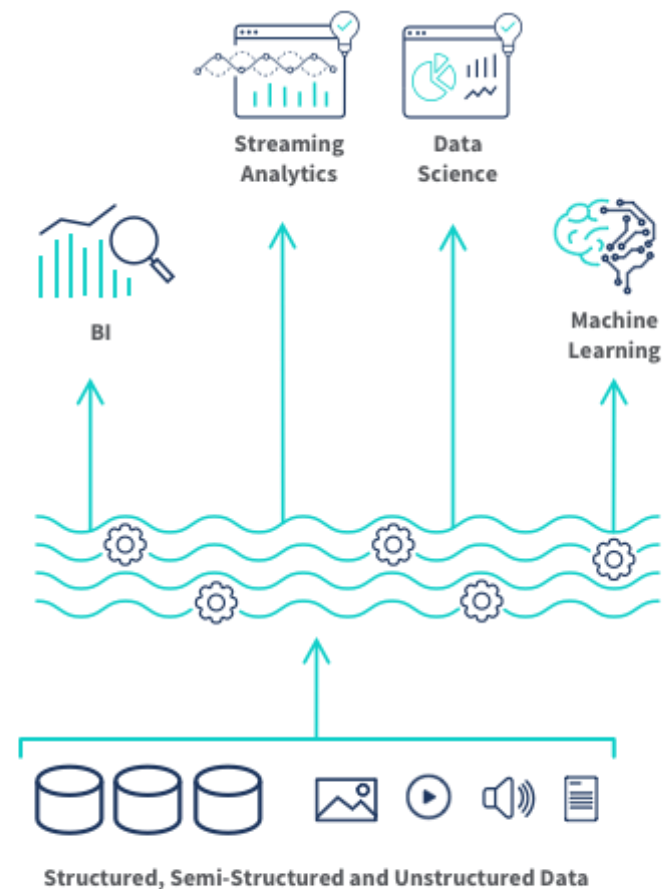
2011

Data Lake



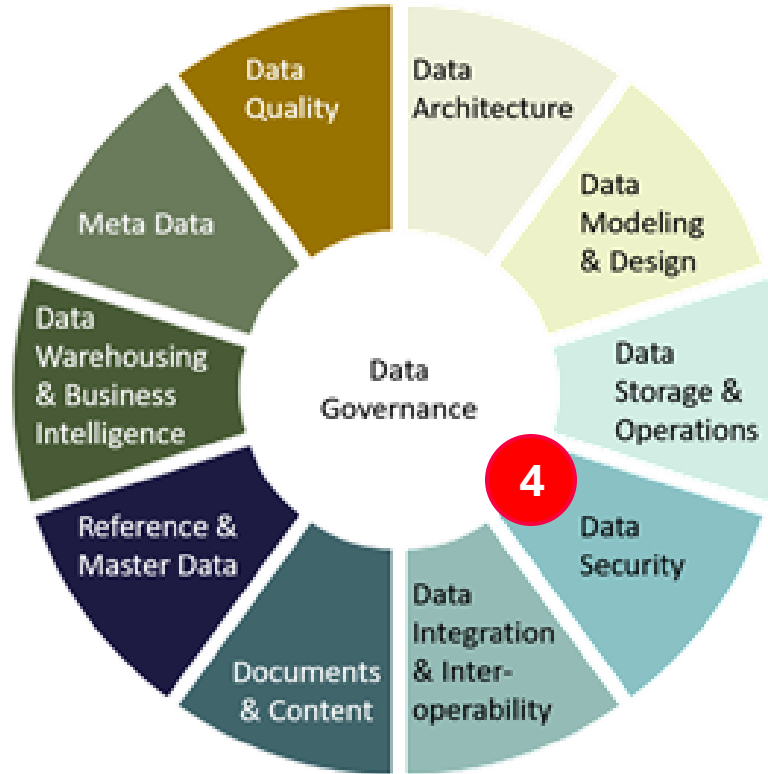
2020

Lakehouse

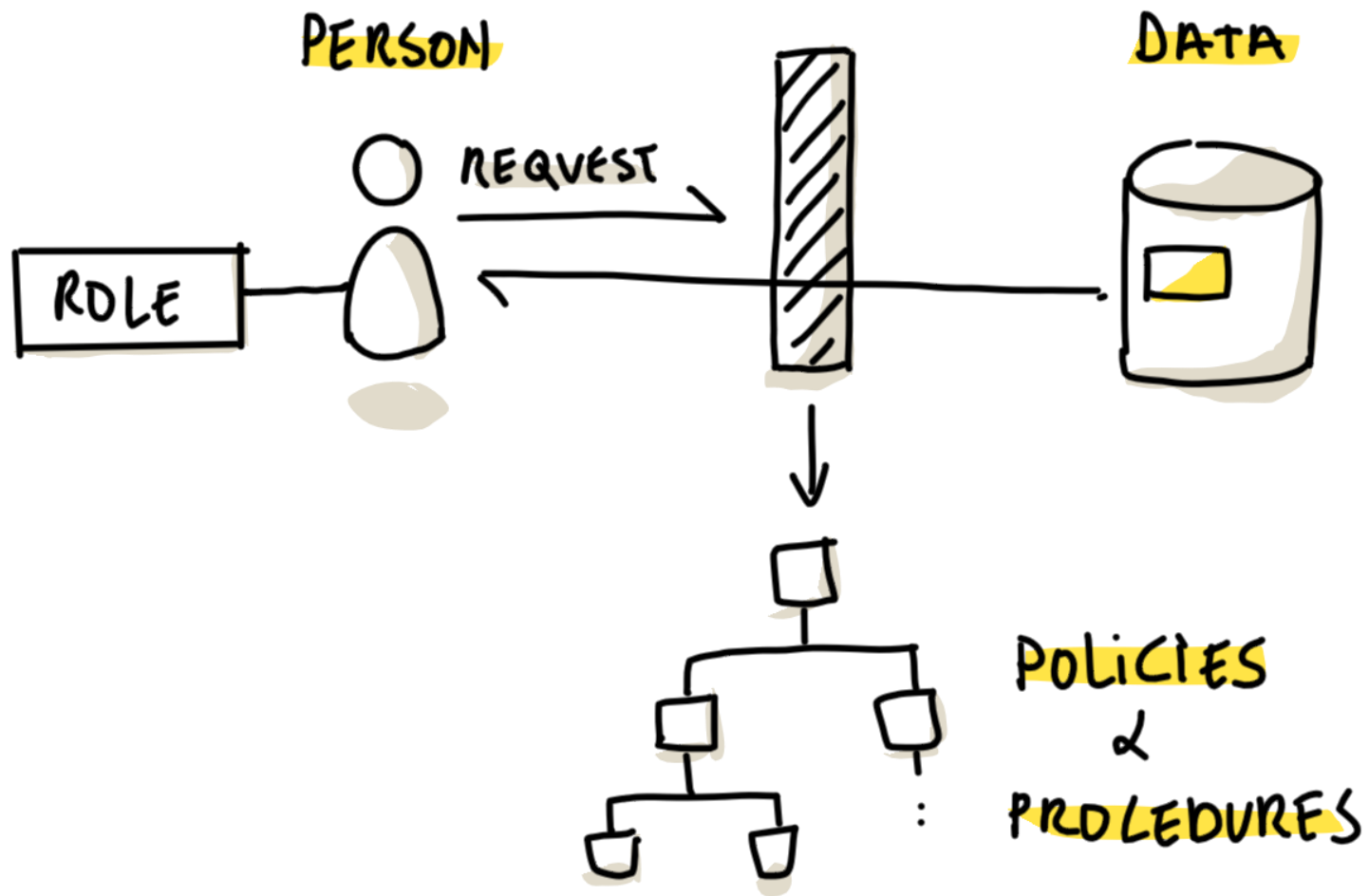


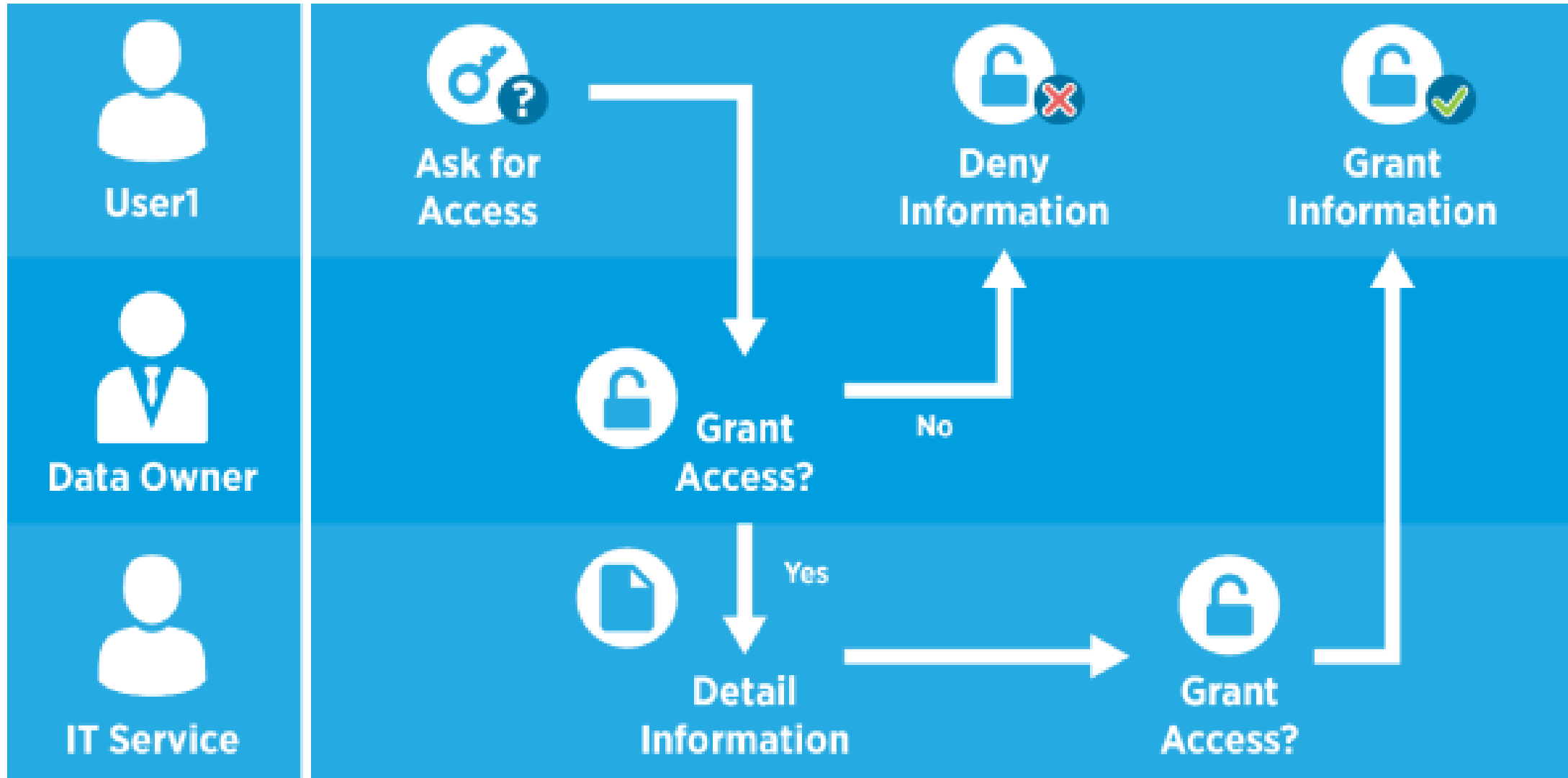
	Database	DWH Database	Data Lake	Data Lakehouse
Cost	+++	+++	+	++
Agility	+	+	++++	+++
Users	Anyone	IT / Business Users	Data Scientists	Anyone
Scaling	Vertical (expensive)	Horizontal (cheaper)		Horizontal (cheaper)
Volume	++	+++	+++++	+++++
Type of Data	Structured	Structured & Semi-Structured	Structured, Semi-Structured & Unstructured	Structured, Semi-Structured & Unstructured
Read Performance	++++ (Depending on the type)	++++	++	+++(+)

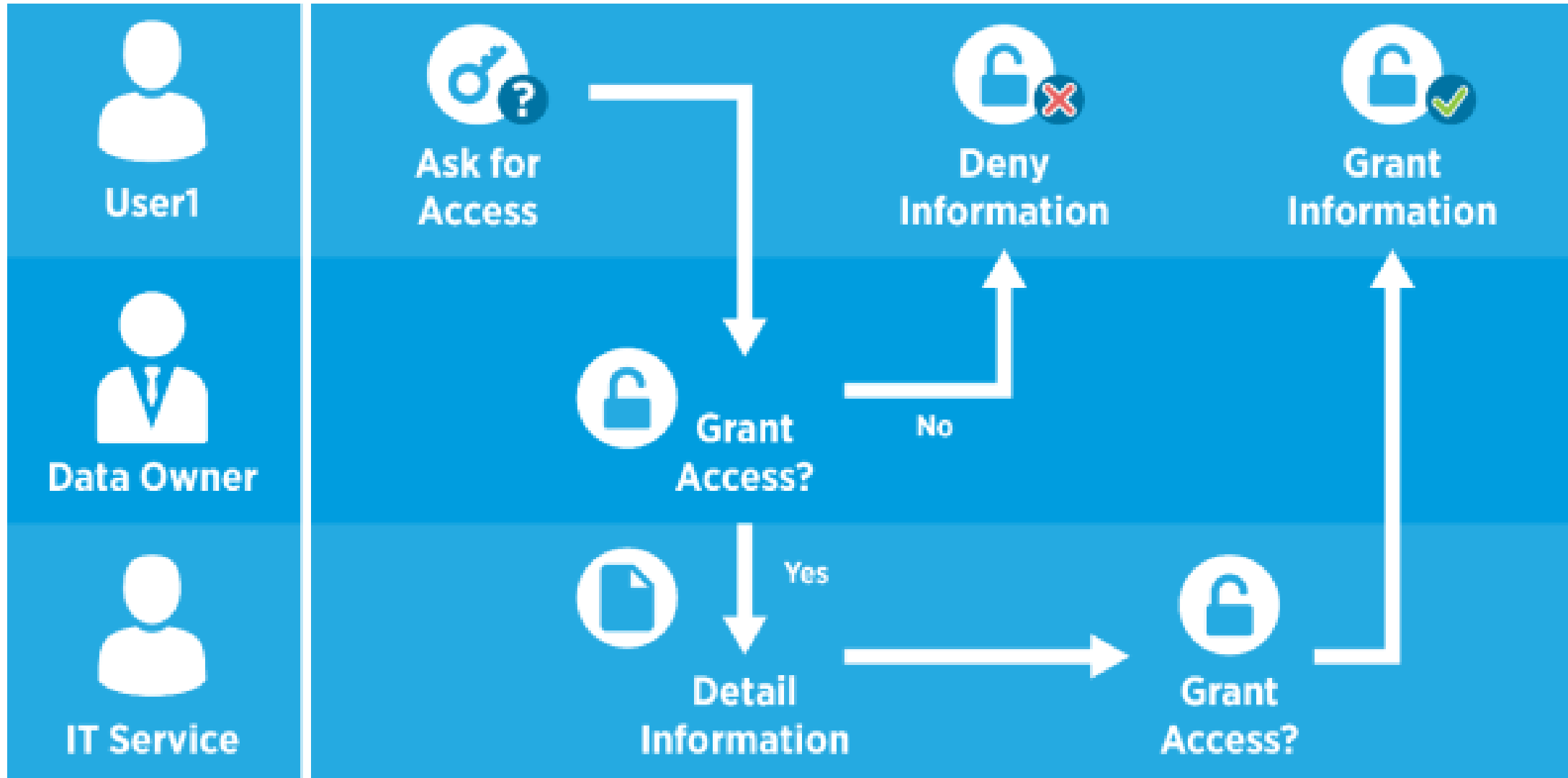
Data Security



Defining Security Policies and
Procedures to provide Data
Access







Multiple Weeks?

Step 3

What are the data objects you need access to?

Data object	Permissions
SALES Schema	Read X ▼
PURCHASING Schema	Read X ▼



Request received



Implementation



Closed



Notifications



Nick Nguyen sent a request that requires your approval #109
3 hours ago

Access control

Owners

Action



DATA_ANALYST
32 data objects



NP

TT

Implement



FINANCE
16 data objects

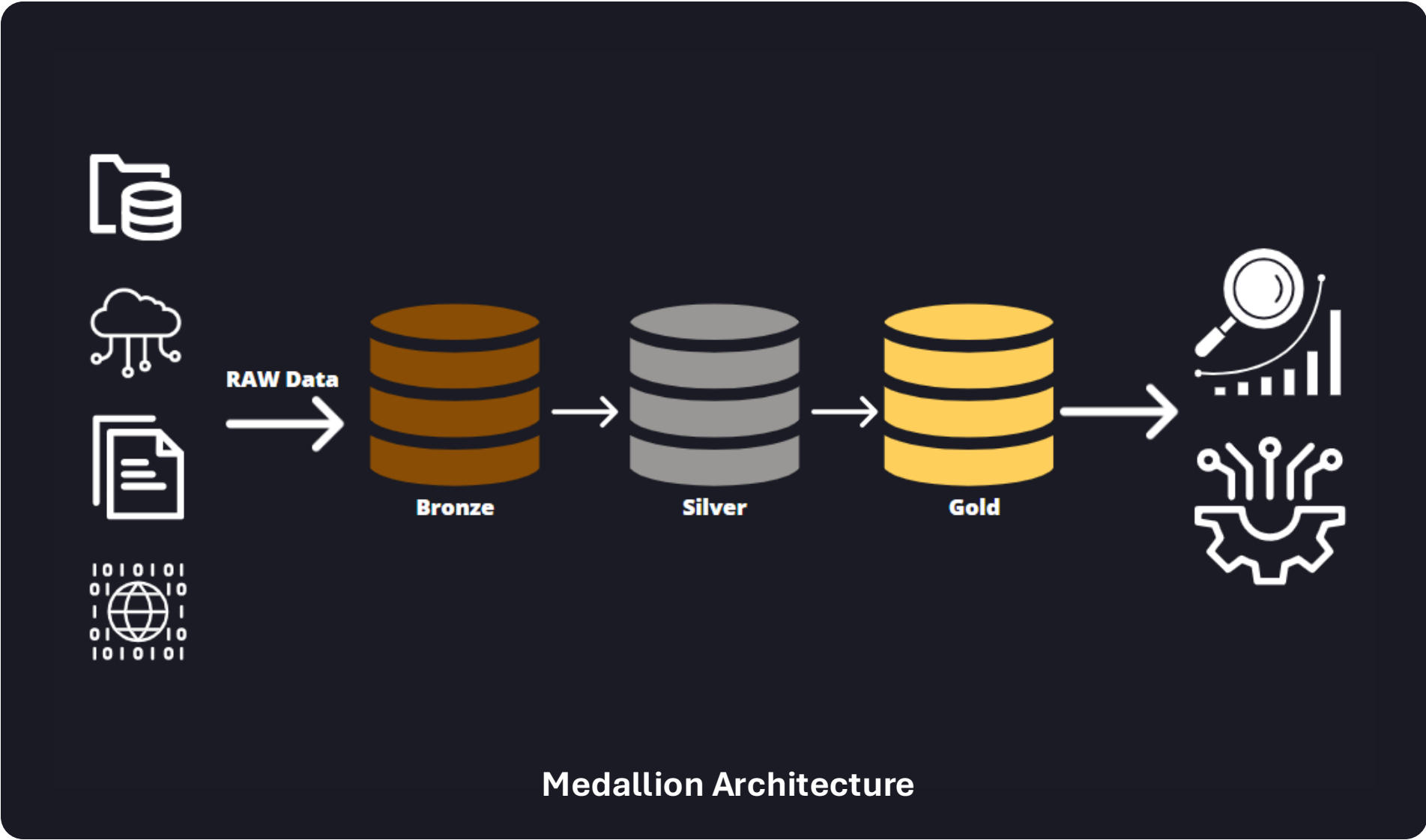


Jin Doe
jin@raito.io

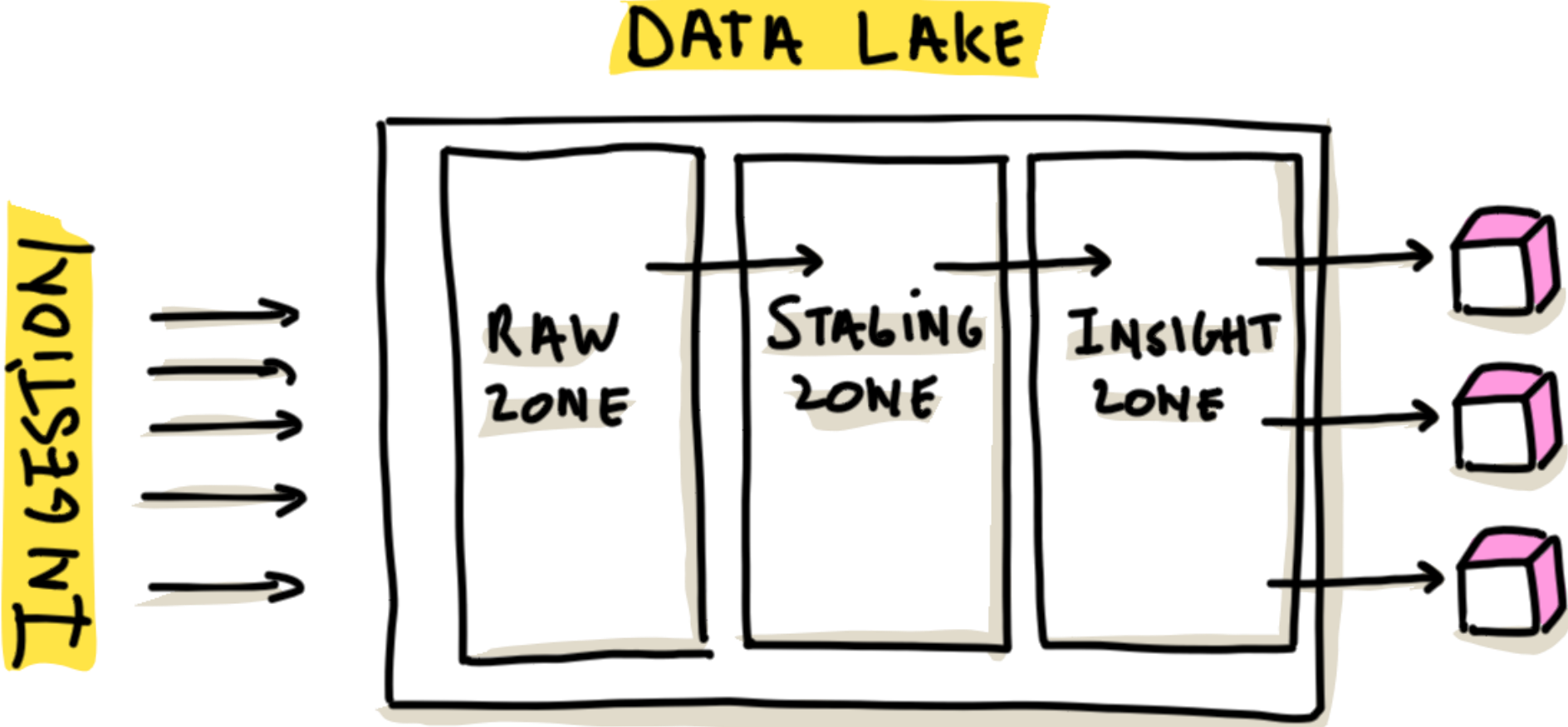
Implement

Policy in a Data Lake

Example Policy



Policy in a Data Lake



Medallion Architecture

Example Data Policy

	RAW ZONE	STAGING ZONE	INSIGHT ZONE
TYPES OF DATA	<ul style="list-style-type: none">• Any kind of data including unstructured data• Examples: videos, tekst files, .csv files, ...	<ul style="list-style-type: none">• Known and structured data• Data from multiple sources is likely to be joined here• Data engineers prep and cleanse data	<ul style="list-style-type: none">• Known, enriched, integrated and cleaned data• Privacy controles like removing personal data
ACCESS	<ul style="list-style-type: none">• Very restricted access• Likely a handful of people or just an admin	<ul style="list-style-type: none">• More access• Mostly data engineers	<ul style="list-style-type: none">• Highest level of access• Most if not all data analysts/scientists

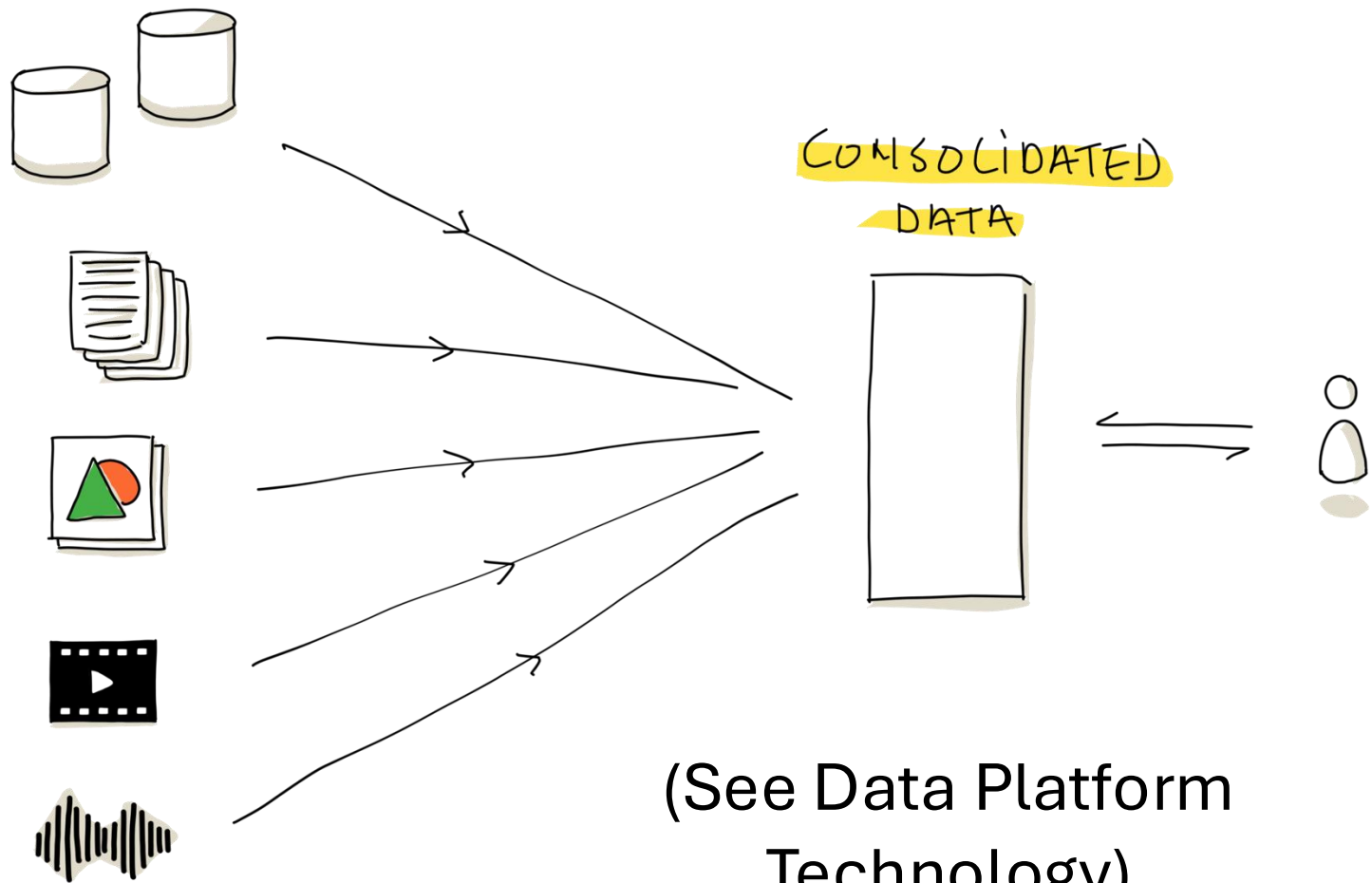
Data Integration & Interoperability



The Movement and Consolidation of Data Within and Between Data Stores, Applications and Organizations

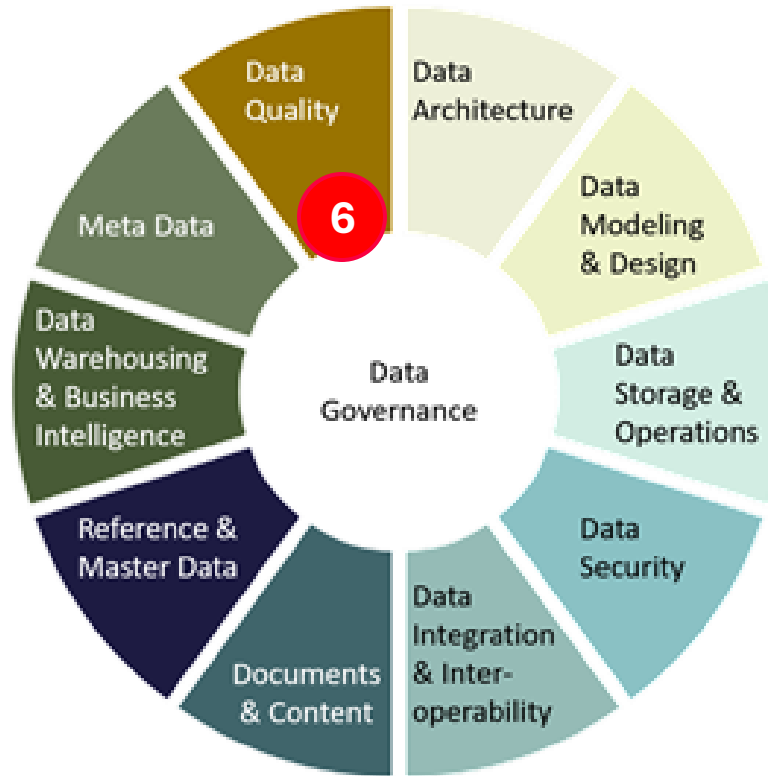


SOURCES



(See Data Platform Technology)

Data Quality



Defining standards and Data Quality controls. Implementing processes to manage and improve DQ



DQ : Example

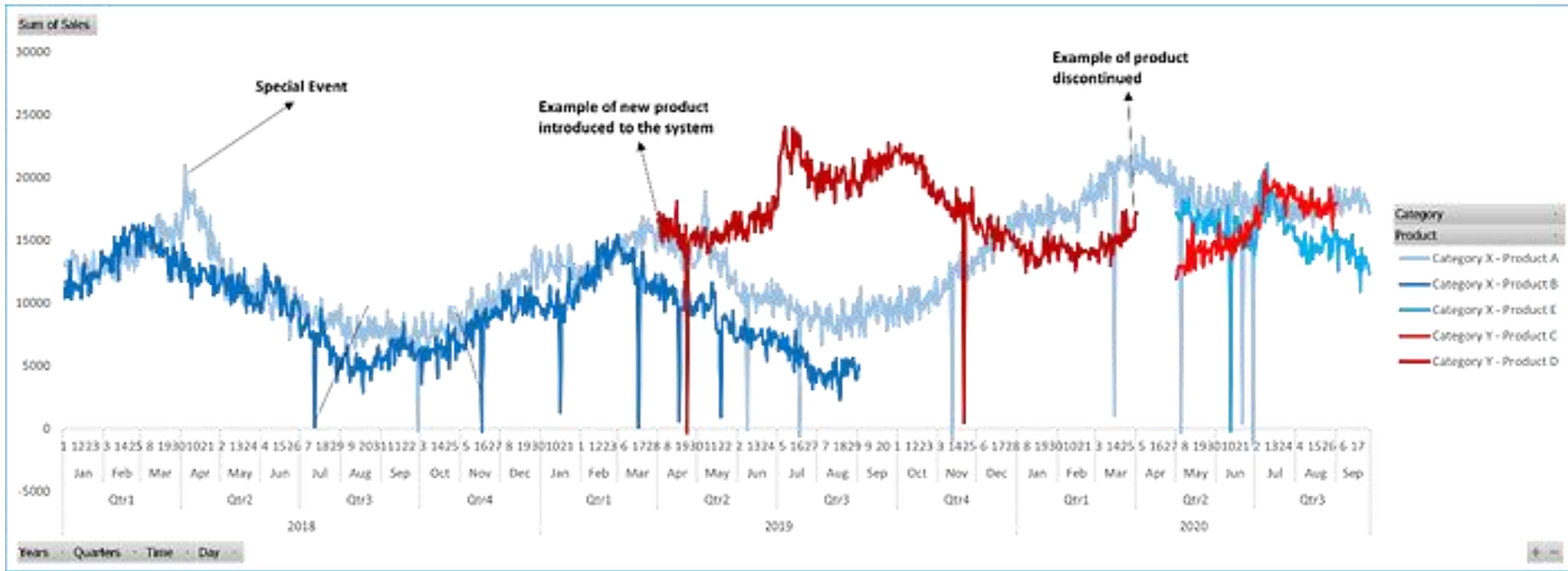
	Sheets	Charts	SmartArt Graphics	WordArt				
	A	B	C	D	E	F	G	H
1	NAME	PHONE	BILLINGSTREET	BILLINGCITY	BILLINGSTATE	WEBSITE		
2	GenePoint	(650) 867-3450	345 Shoreline Park Mountai	Mountain View	CA	www.genepoint.com		
3	United Oil & Gas, Singapore	6504508810	9 Tagore Lane Singapore, S	Singapore	Singapore	http://www.uos.com		
4	Edge Communications	(512) 757-6000	312 Constitution Place Aust	Austin	TX	http://edgecomm.com		
5	Burlington Textiles Corp of America		525-G. Lewis		NC	www.burlington.com		
6	Pyramid Construction Inc.	427-4427	2 Place Juss			www.pyramid.com		
7	Dickenson plc	785-241-6200	1301 Hoch l		KS	dickenson-consulting.com		
8	Grand Hotels & Resorts Ltd	(312) 596-1000	2334 N. Michigan Avenue, S	Chicago	IL	www.grandhotels.com		
9	Express Logistics and Transport	1(503) 421-7800	620 SW 5th Avenue Suite 4	Portland	Oregon	www.expressl&t.net		
10	University of Arizona	77390	888 N Euclid Hallis Center,	Tucson	Arizona			
11	United Oil & Gas	212-8425500	1301 Avenue of the Americ	New York	New York			
12	sForce	ext. 7000	The Landmark @ One Mark	San Francisco	CA			
13								
14								
15								

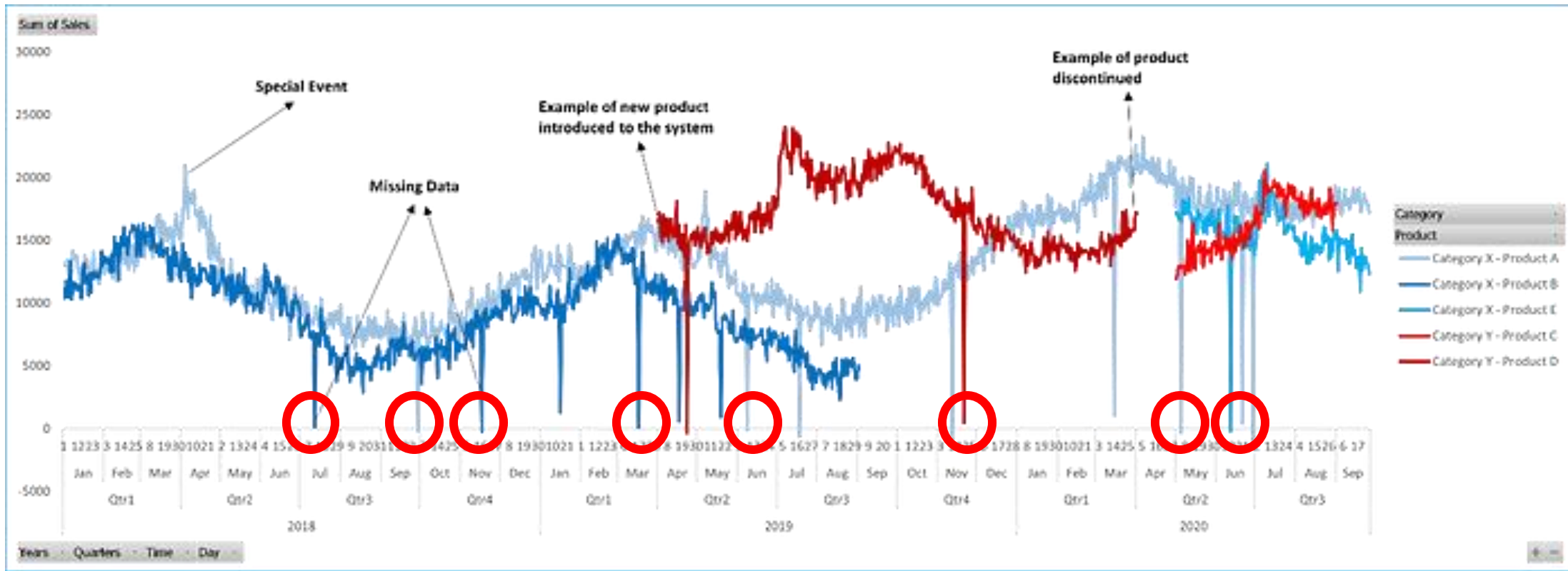
Not Standardized (points to rows 2-4)

Not Complete (points to row 5)

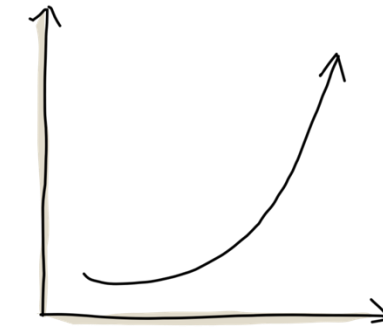
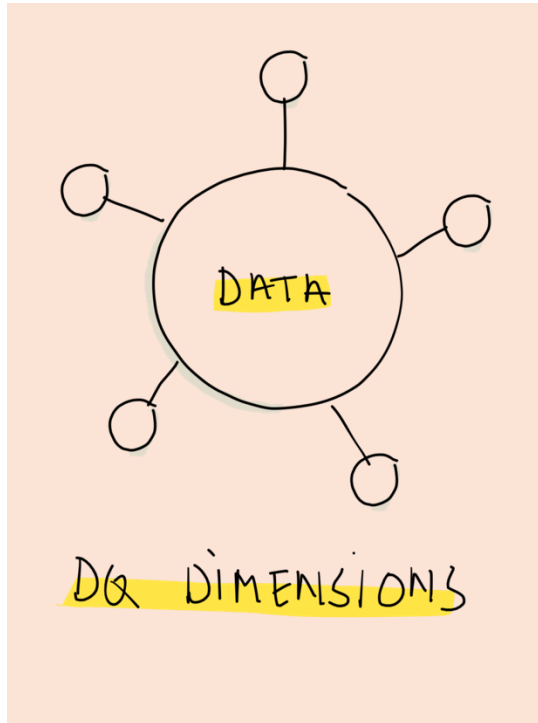
Not Vaild (points to row 12)

Not Consistent (points to row 11)





Data Quality

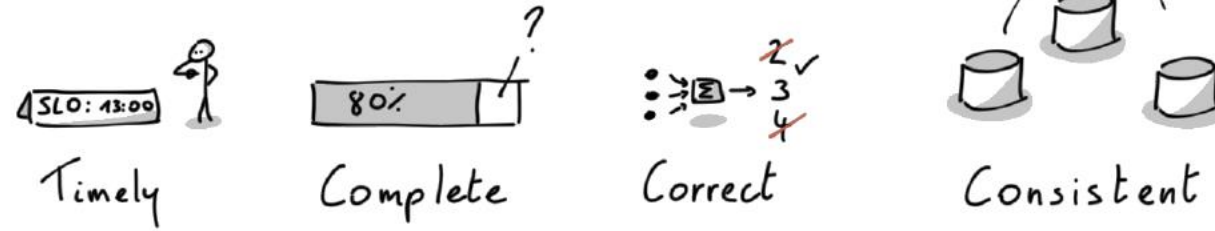


DQ PROCESS
& TOOLS

Introduced in the Previous Lesson



Quality Dimensions (Spotify)



- **Timely:** data is on time (SLO)
- **Complete:** all required data is available
- **Correct:** correct w.r.t. specs, input produces output
- **Consistent:** same meaning across systems, data in sync



What are Data Quality Dimensions?

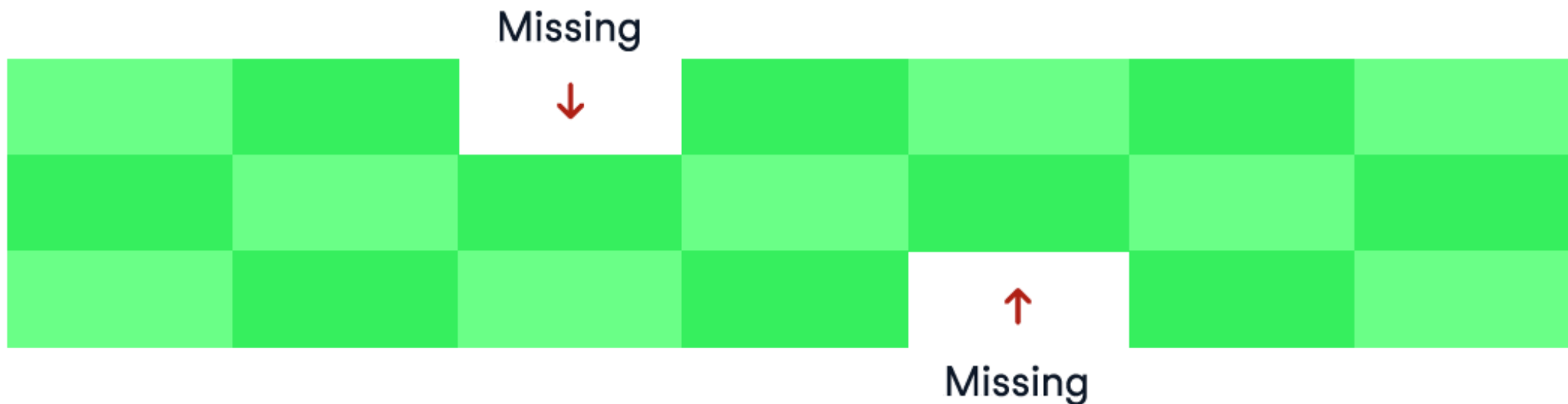
Data Quality is a measurement of the degree to which data is fit for purpose. Good data quality generates trust in data. Data Quality Dimensions are a measurement of a specific attribute of a data's quality.

Completeness – Validity – Accuracy – Uniqueness
Timeliness – Consistency



> Completeness

Completeness measures the degree to which all expected records in a dataset are present. At a data element level, completeness is the degree to which all records have data populated when expected.





Completeness Example

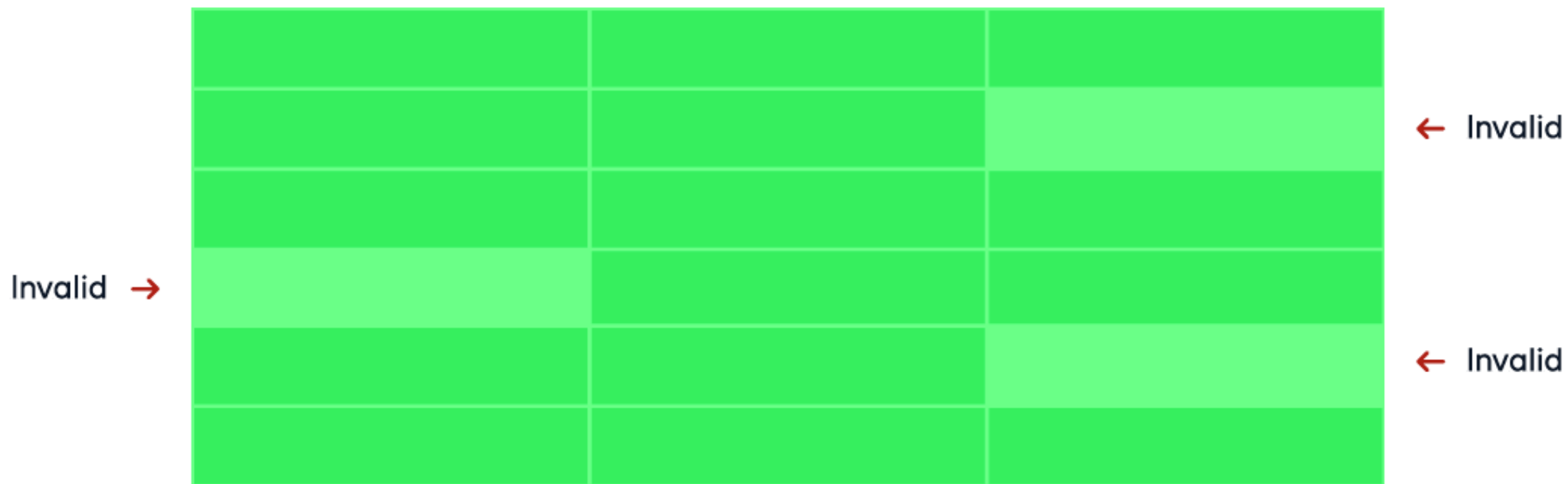
All records must have a value populated in the CustomerName field.

CustomerID	CustomerName	CustomerBirthDate	CustomerAccountType	CustomerAccountBalance	LatestAccountOpenDate
100000192	Robert Brown	4/12/2000	Loan	40390.00	12/20/2026
100000198	Maria Irving	12/1/2025	Deposit	-13280.00	10/21/2018
100000120	Ava Shiffer	10/31/1990	Credit Card	320	3/1/2020
100000192	Robert Brown	4/12/2000	Deposit	40390.00	12/20/2026
100000124	Matthew Martin	5/9/1965	Deposit	70102.00	5/4/2022
100000149		2/4/1988	Loan	0.00	9/20/1990



Validity

Validity measures the degree to which the values in a data element are valid.



Validity Example



- CustomerBirthDate value must be a date in the past.
- CustomerAccountType value must be either Loan or Deposit.
- LatestAccountOpenDate value must be a date in the past.

CustomerID	CustomerName	CustomerBirthDate	CustomerAccountType	CustomerAccountBalance	LatestAccountOpenDate
100000192	Robert Brown	4/12/2000	Loan	40390.00	12/20/2026
100000198	Maria Irving	12/1/2025	Deposit	-13280.00	10/21/2018
100000120	Ava Shiffer	10/31/1990	Credit Card	320	3/1/2020
100000192	Robert Brown	4/12/2000	Deposit	40390.00	12/20/2026
100000124	Matthew Martin	5/9/1965	Deposit	70102.00	5/4/2022
100000149		2/4/1988	Loan	0.00	9/20/1990



Accuracy

Accuracy measures the degree to which data is correct and represents the truth.

Verified Source Document

Orange	Orange	Orange
Green	Green	Green
Blue	Blue	Blue
Purple	Purple	Purple

Downstream Table

Orange	Orange	Orange
Green	Green X	Green
Blue	Blue	Blue
Purple	Purple	Purple



Accuracy Example

All records in the Customer Table must have accurate Customer Name, Customer Birthdate, and Customer Address fields when compared to the Tax Form.

Tax Form

Name: Ava Shiffer Birthdate: 10/30/1990

Address: 910 Quality St

City: Washington State: DC

Zip: 20008



CustomerName	CustomerBirthDate	CustomerAddress	CustomerCity	CustomerState	CustomerZip
Ava Shiffer	10/31/1990	910 Quality St	Washington	WA	20008



After I sent **a late notice about an outstanding invoice** to a third-party firm I sub-contract for, we discovered that while the check was indeed in the mail, unfortunately it was mailed to the wrong address—a valid but inaccurate address.



Uniqueness

Uniqueness measures the degree to which the records in a dataset are not duplicated.





Uniqueness Example

All records must have a unique CustomerID and CustomerName.

CustomerID	CustomerName	CustomerBirthDate	CustomerAccountType	CustomerAccountBalance	LatestAccountOpenDate
100000192	Robert Brown	4/12/2000	Loan	40390.00	12/20/2026
100000198	Maria Irving	12/1/2025	Deposit	-13280.00	10/21/2018
100000120	Ava Shiffer	10/31/1990	Credit Card	320	3/1/2020
100000192	Robert Brown	4/12/2000	Deposit	40390.00	12/20/2026
100000124	Matthew Martin	5/9/1965	Deposit	70102.00	5/4/2022
100000149		2/4/1988	Loan	0.00	9/20/1990



Timeliness

Timeliness is the degree to which a dataset is available when expected and depends on service level agreements being set up between technical and business resources.

SLA	Table Load Time
08:00 am	07:59 am
10:00 am	09:59 am
11:00 am	11:01 am

← Missed the SLA



Timeliness Example

All records in the customer dataset must be loaded by the 9:00 am.

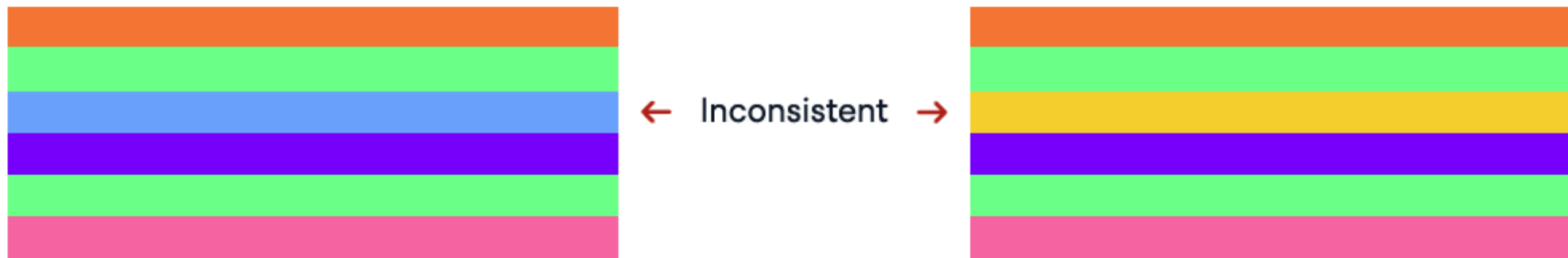


CustomerID	CustomerName
100000192	01-01-2023 11:07 am
100000198	01-01-2023 11:07 am
100000120	01-01-2023 11:07 am



Consistency

Consistency is a data quality dimension that measures the degree to which data is the same across all instances of the data. Consistency can be measured by setting a threshold for how much difference there can be between two datasets.

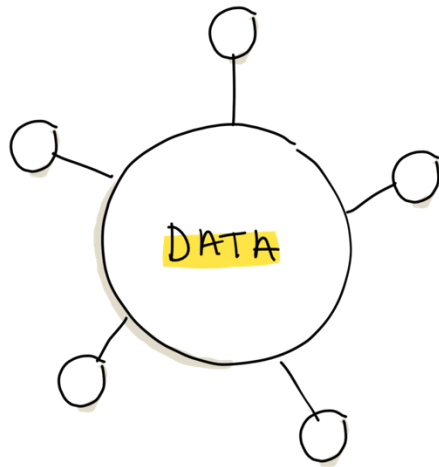




Consistency Example

AccountTableCustomerID	CustomerTableCustomerID
108394858	108394858
192039482	192039482
203475849	NULL X
2930485953	NULL X
102832748	102832748

Data Quality



DQ DIMENSIONS

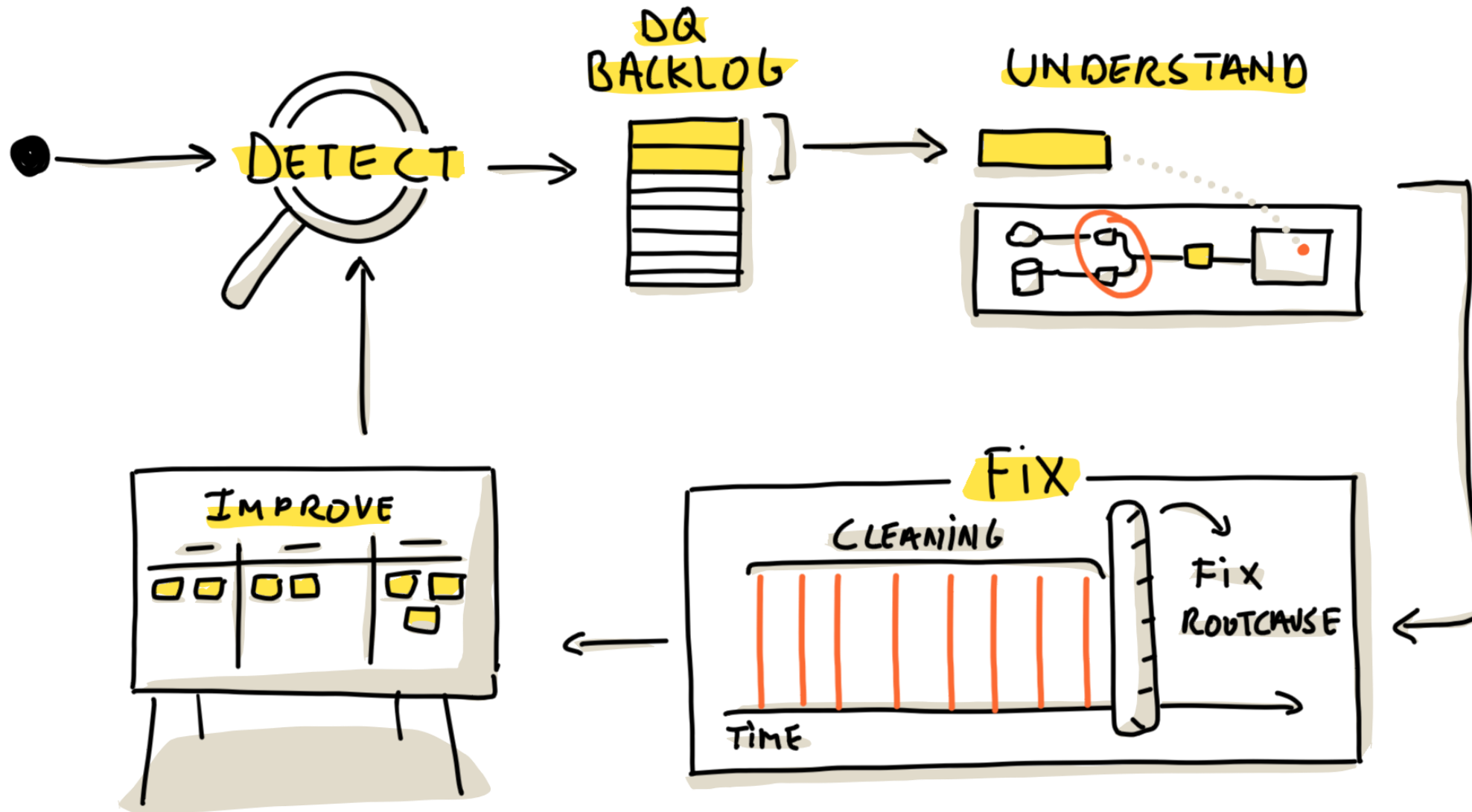


DQ PROCESS

& TOOLS

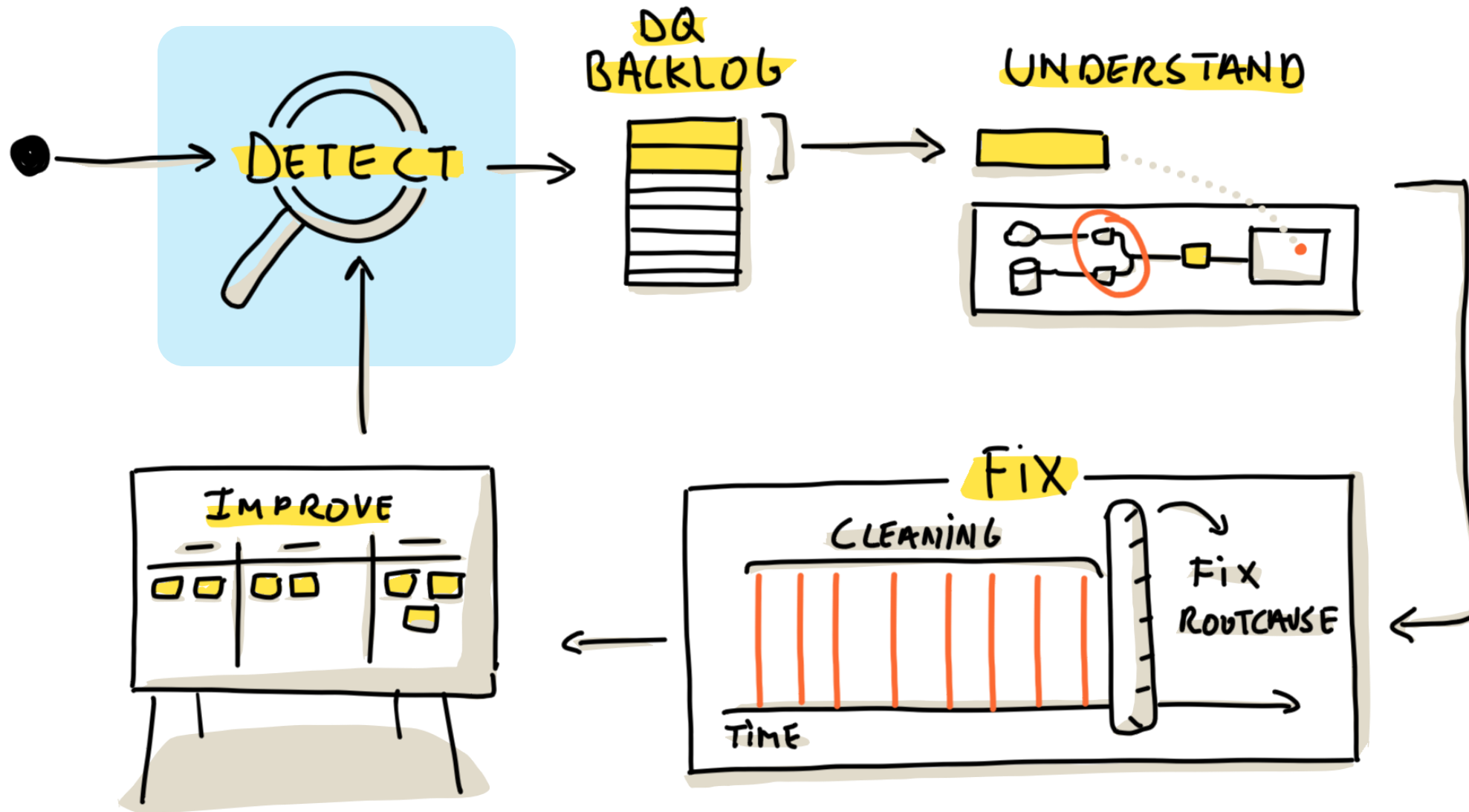


Data Quality Process





Data Quality Process





DQ Tests

SQL customer_id.assert.sql

```
1 -- check if customers contains null values
2 SELECT customer_id as customer_id,
3        customer_name as customer_name
4 FROM customers
4 WHERE customer_id IS NULL
```

run

if 0 row returned:

Assertion passed ✓

if >=1 row(s) returned:

Assertion failed ✗



DQ Monitoring



Home > Knowledge Catalog > Sources > MDM

party_full

Use In

Overview Profile Data Quality Data Preview Lineage Relationships 2999 Records 7 Attributes Profiled 2 mins

Filter attributes, values, masks

Name	Terms	Insights	Top 3 Values	Mask Analysis
<u>src_primary_key</u>		3 Duplicates	3% NNN 0% 145 0% 146	3% LLL 47% DDD 50% DDDD
<u>src_name</u>	<u>Last Name</u>	3 Duplicates	24% Null 3% Green 2% Kazmer	6% LLLL 5% LLLLL Show All +29
<u>src_sin</u>	<u>Social Insurance Number</u>	NULL 24%	24% Null 0% 103792776 0% SIN: 999670052	24% LLL: DDDDDDDI 18% DDDDDDDDD Show All +22
<u>src_card</u>	<u>Credit Card Number</u>	7 Exceptions	2% ##### 1% ##### 0% #####	98% DDDDDDDDDDDDI 2% LLLL



Home > Knowledge Catalog

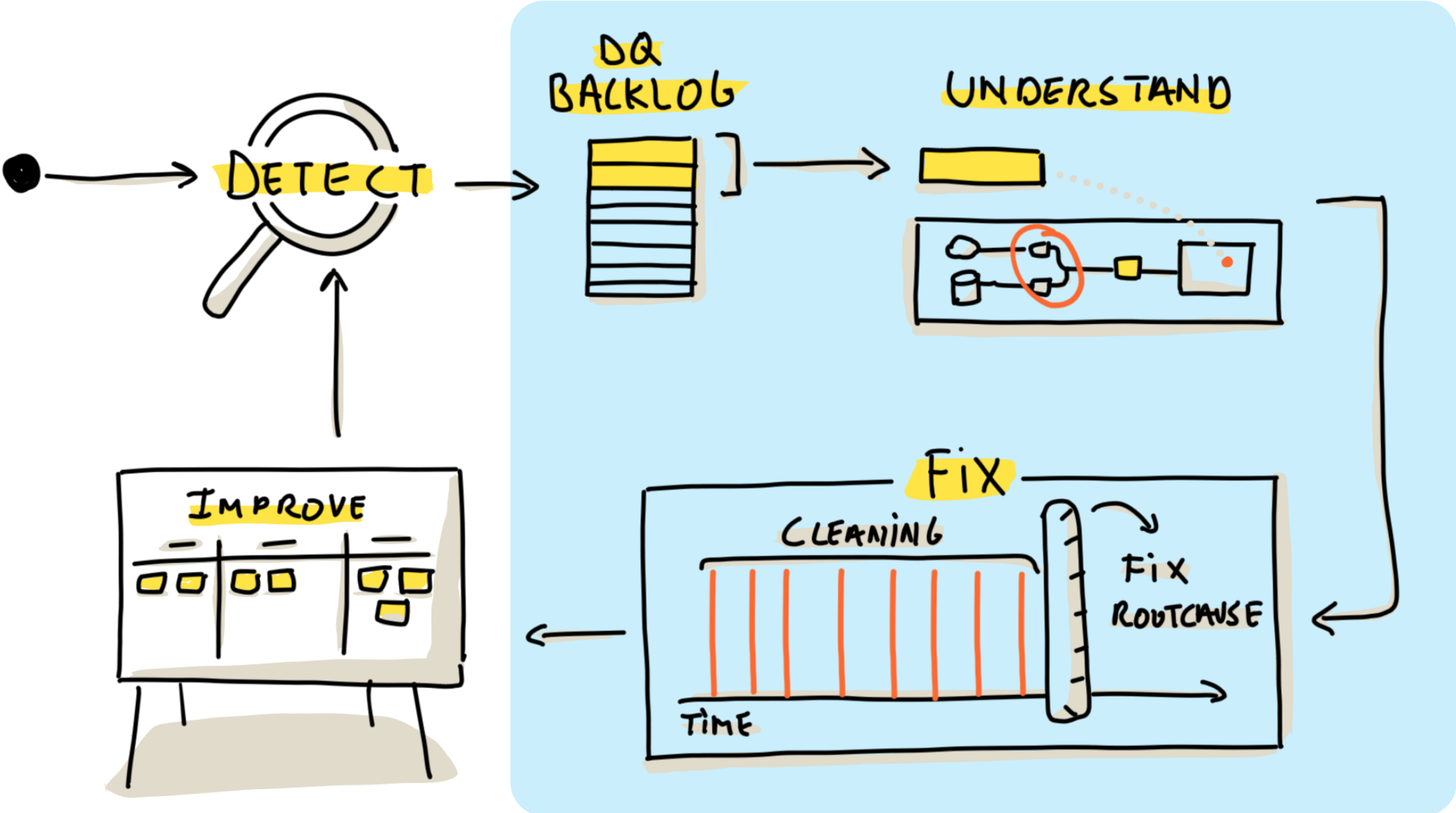
Data Assets

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 style="width: 50%; background-color: green;"></div> <div style="width: 50%; background-color: red;"></div>	
<input type="checkbox"/>	<u>Master customer</u>	PII Customer	<div style="width: 20%; background-color: green;"></div> <div style="width: 80%; background-color: red;"></div>	
<input type="checkbox"/>	<u>Customers 2019</u>	PII Customer	<div style="width: 95%; background-color: green;"></div> <div style="width: 5%; background-color: red;"></div>	
<input type="checkbox"/>	<u>comp</u>	Account	<div style="width: 80%; background-color: green;"></div> <div style="width: 20%; background-color: red;"></div>	
<input type="checkbox"/>	<u>Customer campaigns</u>	Customer Campaign	<div style="width: 95%; background-color: green;"></div> <div style="width: 5%; background-color: red;"></div>	
<input type="checkbox"/>	<u>cstmr</u>	PII Customer	<div style="width: 95%; background-color: green;"></div> <div style="width: 5%; background-color: red;"></div>	
<input type="checkbox"/>	<u>employees_2020</u>	PII Employee	<div style="width: 50%; background-color: green;"></div> <div style="width: 50%; background-color: red;"></div>	
<input type="checkbox"/>	<u>Master address</u>	Address	<div style="width: 20%; background-color: green;"></div> <div style="width: 80%; background-color: red;"></div>	
<input type="checkbox"/>	<u>cstomers_2019_ext</u>	PII Customer	<div style="width: 95%; background-color: green;"></div> <div style="width: 5%; background-color: red;"></div>	
<input type="checkbox"/>	<u>account_list</u>	PII Account	<div style="width: 80%; background-color: green;"></div> <div style="width: 20%; background-color: red;"></div>	



Data Quality Process

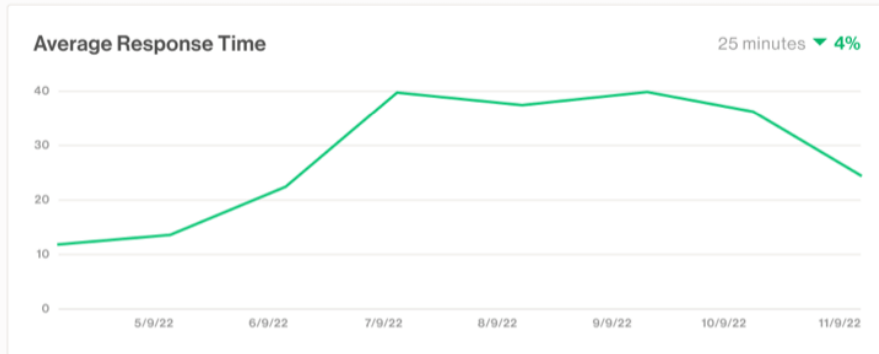
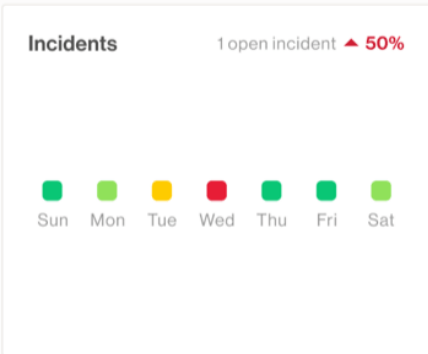
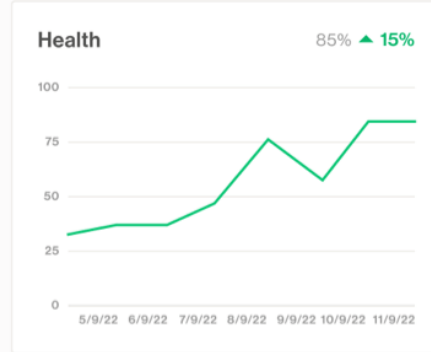
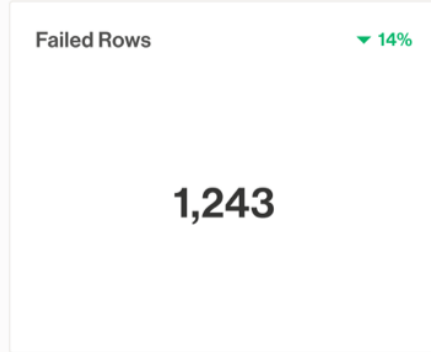
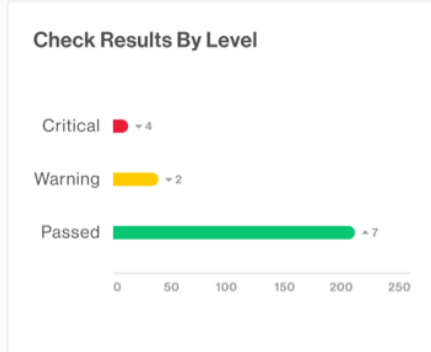
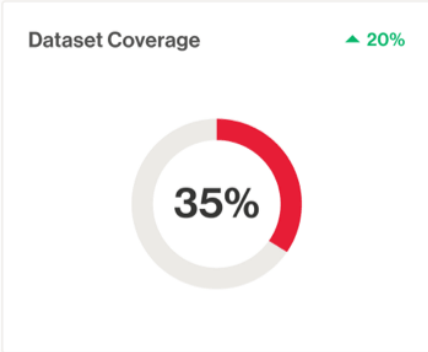




Dashboard

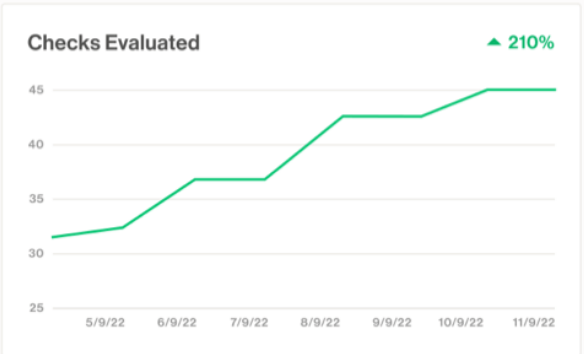
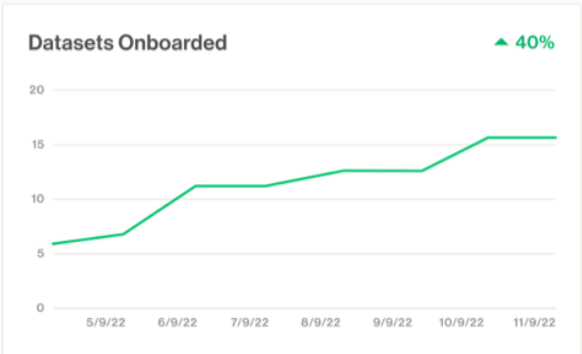
Last 7 days

Finance



Incident Resolution Leaderboard

1	JD John Doe	13
2	KM Kelly Madison	12
3	JF Jen Finley	7
4	TD Thomas Davidson	4
5	MT Marc Tyler	3



Meta Data



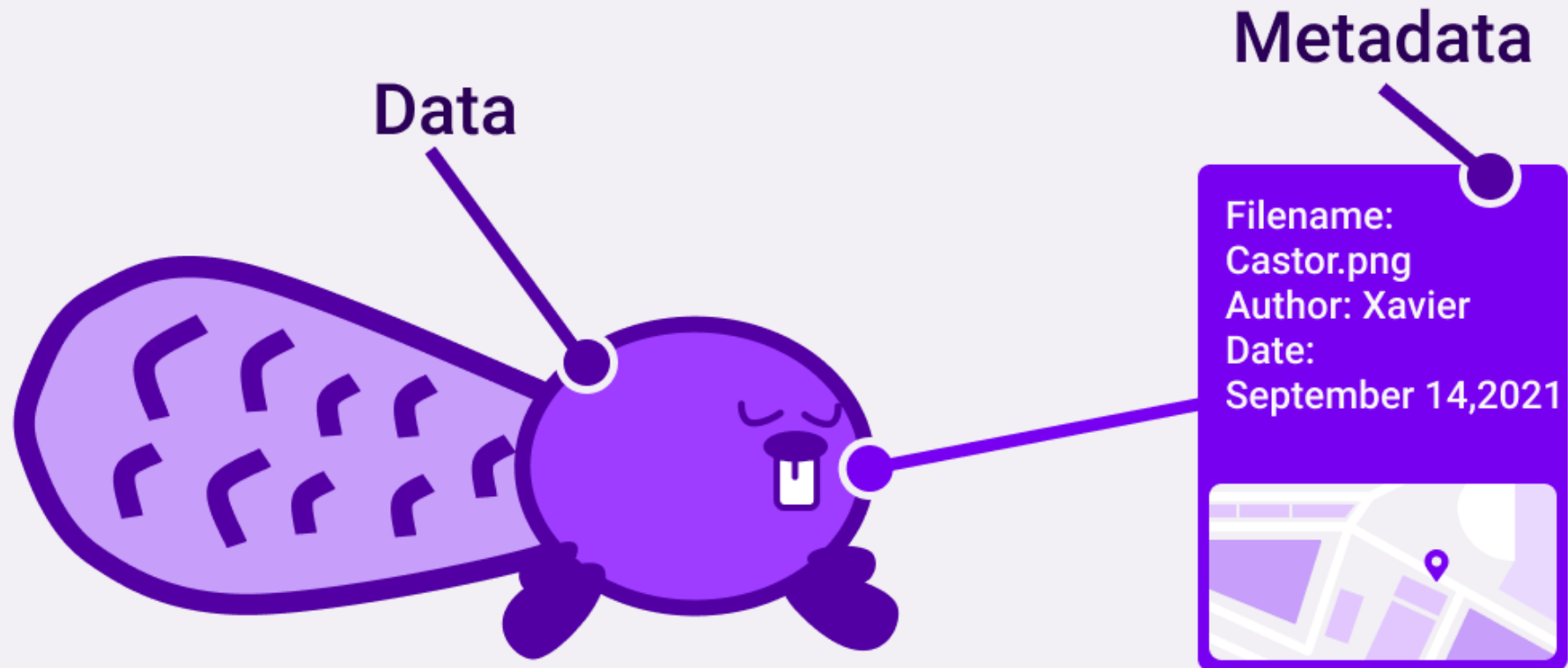
...



META-DATA PORTAL



Metadata = “Data about Data”



An iceberg floating in the ocean. The tip of the iceberg is above the water surface, while the much larger, jagged base is submerged underwater. The sky is blue with some clouds, and the water is a deep blue. The horizon line is visible in the distance.

Data

(What you can use)

Metadata

(Describing the data format, the actions leading to this data, documentation, business context, ...)



Types of Metadata

Technical

Definitional

Schemas, data types,
models, etc.

Operational

Descriptive

How is data being
produced?
(Source /
Transformations /
...)

Business

Descriptive

Data tags,
classifications,
mappings to
business
relationships, etc.

Social

Descriptive

Metadata about
user-generated
content, business
knowledge, etc.



Types of Metadata

Technical

Definitional

Schemas, data types,
models, etc.

Operational

Descriptive

How is data being
produced?
(Source /
Transformations /
...)

Business

Descriptive

Data tags,
classifications,
mappings to
business
relationships, etc.

Social

Descriptive

Metadata about
user-generated
content, business
knowledge, etc.



Table Structure

employee_id	first_name	last_name	nin	department_id
44	Simon	Martinez	HH 45 09 73 D	1
45	Thomas	Goldstein	SA 75 35 42 B	2
46	Eugene	Comelsen	NE 22 63 82	2
47	Andrew	Petculescu	XY 29 87 61 A	1
48	Ruth	Stadick	MA 12 89 36 A	15
49	Barry	Scardelis	AT 20 73 18	2
50	Sidney	Hunter	HW 12 94 21 C	6
51	Jeffrey	Evans	LX 13 26 39 B	6
52	Doris	Bemdt	YA 49 88 11 A	3
53	Diane	Eaton	BE 08 74 68 A	1
54	Bonnie	Hall	WW 53 77 68 A	15
55	Taylor	Li	ZE 55 22 80 B	1

Data

Metadata

Column	Data Type	Description
employee_id	int	Primary key of a table
first_name	nvarchar(50)	Employee first name
last_name	nvarchar(50)	Employee last name
nin	nvarchar(15)	National Identification Number
position	nvarchar(50)	Current position title, e.g. Secretary
department_id	int	Employee department. Ref: Departments
gender	char(1)	M = Male, F = Female, Null = unknown
employment_start_date	date	Start date of employment in organization.
employment_end_date	date	Employment end date. Null if employee sti



Types of Metadata

Technical

Definitional

Schemas, data types,
models, etc.

Operational

Descriptive

How is data being
produced?
(Source /
Transformations /
...)

Business

Descriptive

Data tags,
classifications,
mappings to
business
relationships, etc.

Social

Descriptive

Metadata about
user-generated
content, business
knowledge, etc.



Data Lineage

- Documentation of the flow and transformation of data as it moves from source to product.
- It is crucial for understanding a data product's origins, transformations
- Four key points to define data lineage:
 1. **Data Source Identification:** identifying the original sources of data
 2. **Data Movement Tracking:** It traces how data is transferred and transformed
 3. **Dependency Mapping:** Data lineage maps dependencies between different data elements
 4. **End-to-End Visibility:** It provides a comprehensive view of data's journey



Example: E2E Visible Data Lineage

Tool: DBT (Data Build Tool)

Data Movement Tracking



Data Source:

Product:



Dependency: stg_eltool__orders depends on warehouse.orders



Types of Metadata

Technical

Definitional

Schemas, data types,
models, etc.

Operational

Descriptive

How is data being
produced?
(Source /
Transformations /
...)

Business

Descriptive

Data tags,
classifications,
mappings to
business
relationships, etc.

Social

Descriptive

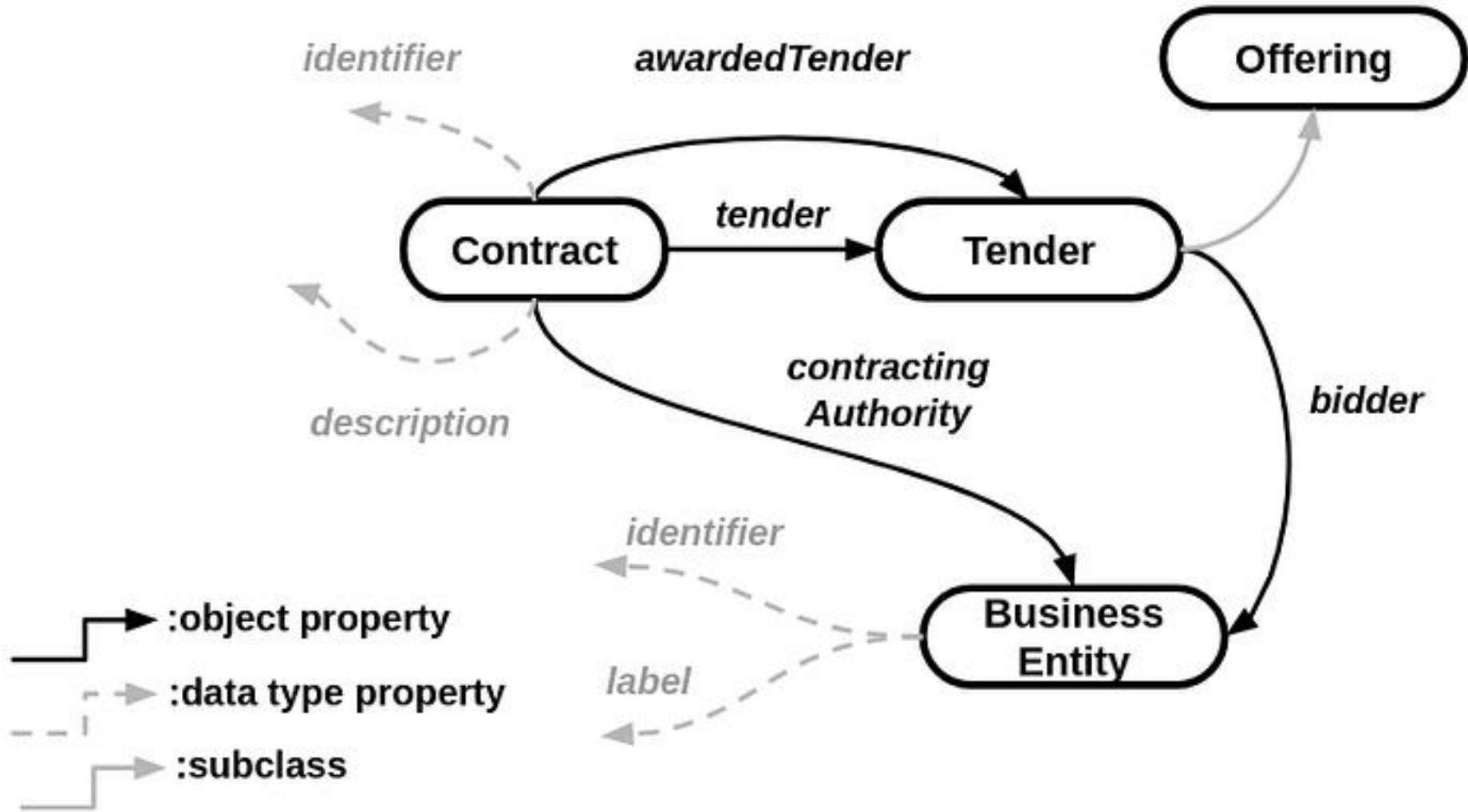
Metadata about
user-generated
content, business
knowledge, etc.



Business Glossary

The screenshot displays the Atlan Business Glossary interface. On the left is a navigation sidebar with sections for 'Assets', 'Glossary', and 'Insights'. The 'Glossary' section is active, showing a search bar and a list of glossaries including 'Aisle', 'Concepts', 'Consumer Product Goods', 'COVID-19', 'Example Glossary', 'Instacart', 'KPIs', and 'Metrics'. Under 'Metrics', 'Customer Acquisition Cost' is selected. The main content area shows the 'Customer Acquisition Cost' term page, which includes a 'Readme' section with two methods: 'Simple method' and 'Complex method'. The 'Simple method' section contains the formula $CAC = \frac{MCC}{CA}$ and a list of definitions for CAC, MCC, and CA. The 'Complex method' section contains the formula $CAC = \frac{MCC + W + S + PS + O}{CA}$ and a list of definitions for CAC, MCC, W, S, and O. On the right side of the page, there is a metadata panel with sections for 'Overview', 'Owners', 'Classification', 'Certificate', 'Categories', 'Related Terms', and 'Custom Metadata'. The 'Owners' section is highlighted with a red box and shows 'chris' as the owner. The 'Certificate' section shows 'Verified' with a green checkmark and 'chris 3 months ago' as the certifier. The 'Related Terms' section shows 'Average Selling Price', 'Churn Rate', and 'Customer Lifetime Value'.

Data Models





Types of Metadata

Technical

Definitional

Schemas, data types,
models, etc.

Operational

Descriptive

How is data being
produced?
(Source /
Transformations /
...)

Business

Descriptive

Data tags,
classifications,
mappings to
business
relationships, etc.

Social

Descriptive

Metadata about
user-generated
content, business
knowledge, etc.

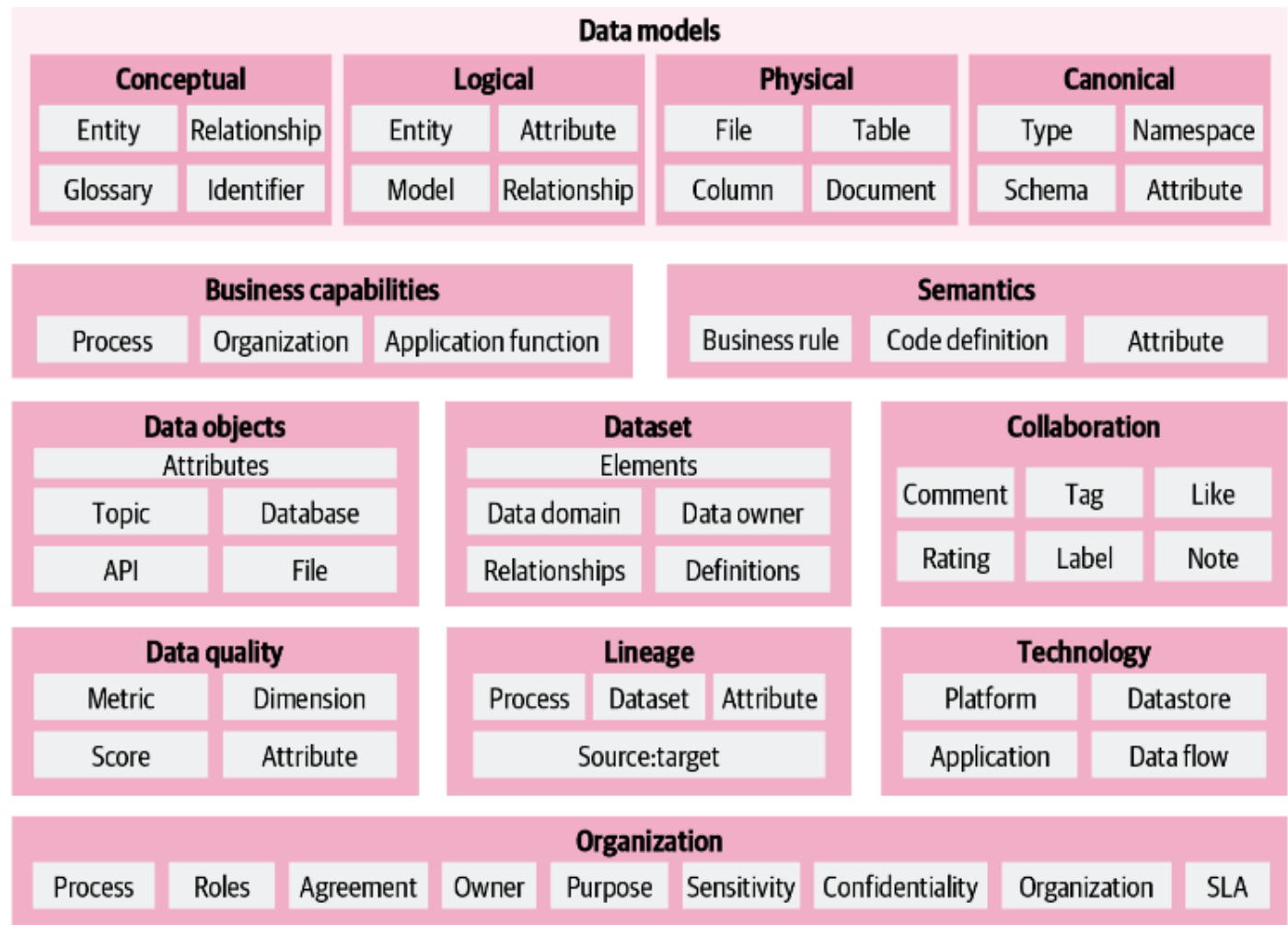
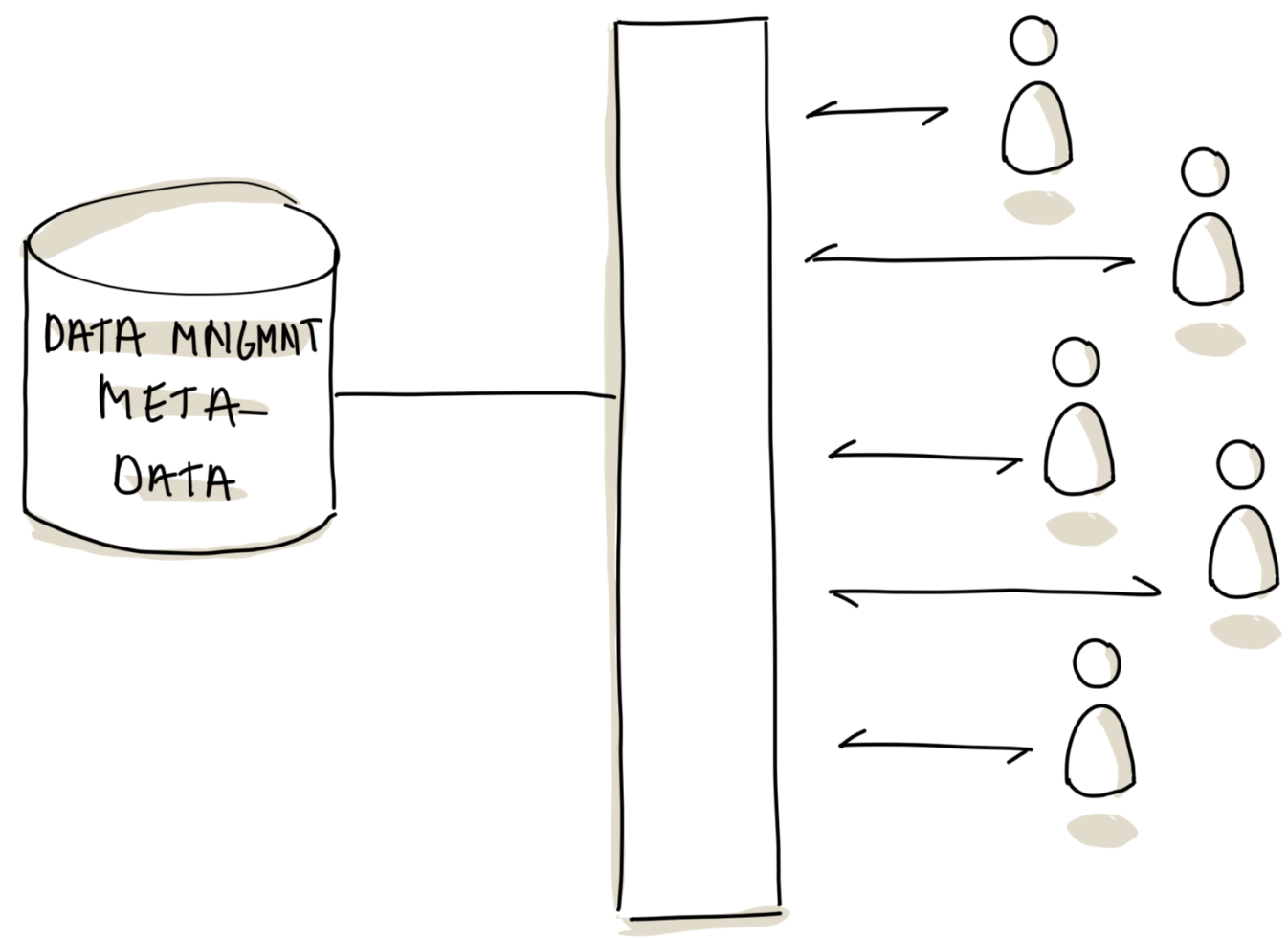


Figure 10-2. An overview showing all the major metadata management subject areas for the enterprise. These subject areas are nonexhaustive. Each organization uses the areas differently.



META-DATA PORTAL





META-DATA PORTAL

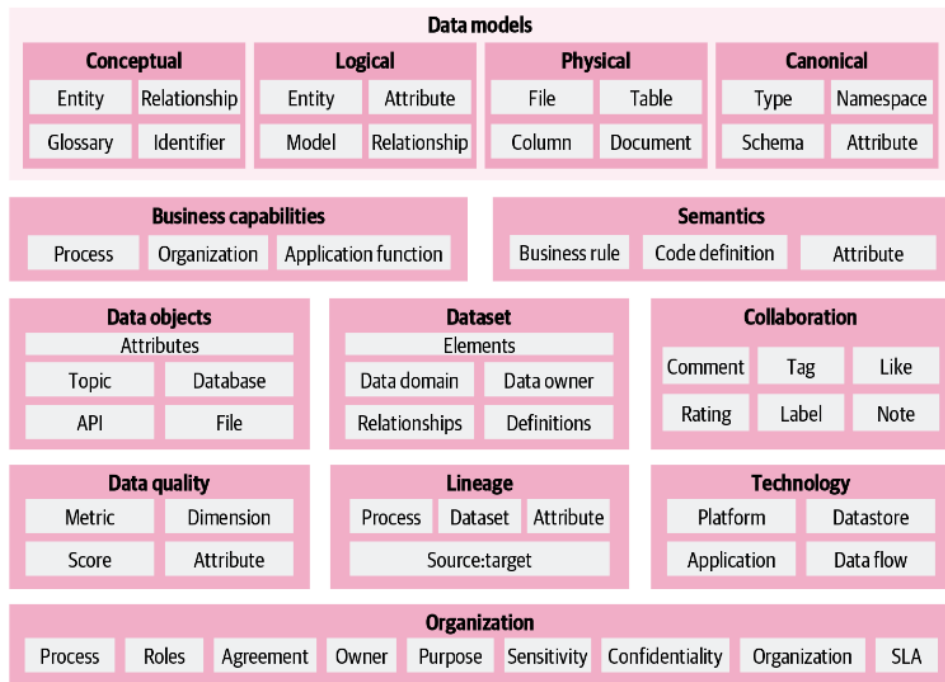
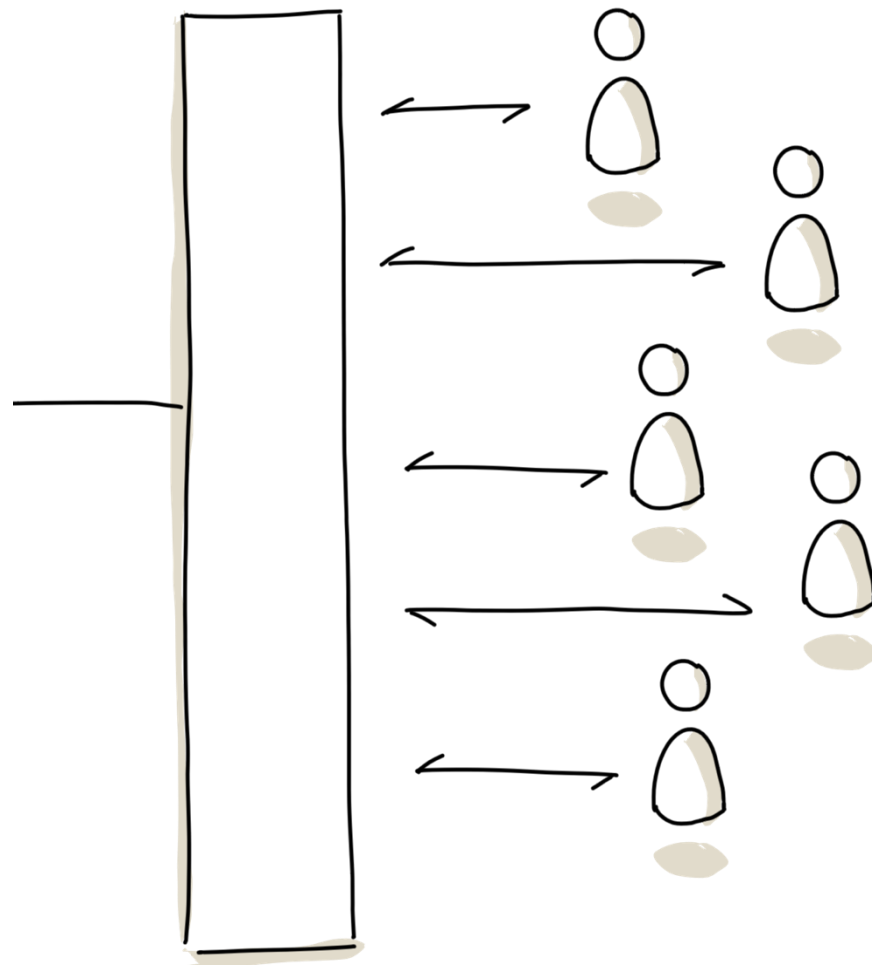


Figure 10-2. An overview showing all the major metadata management subject areas for the enterprise. These subject areas are nonexhaustive. Each organization uses the areas differently.



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

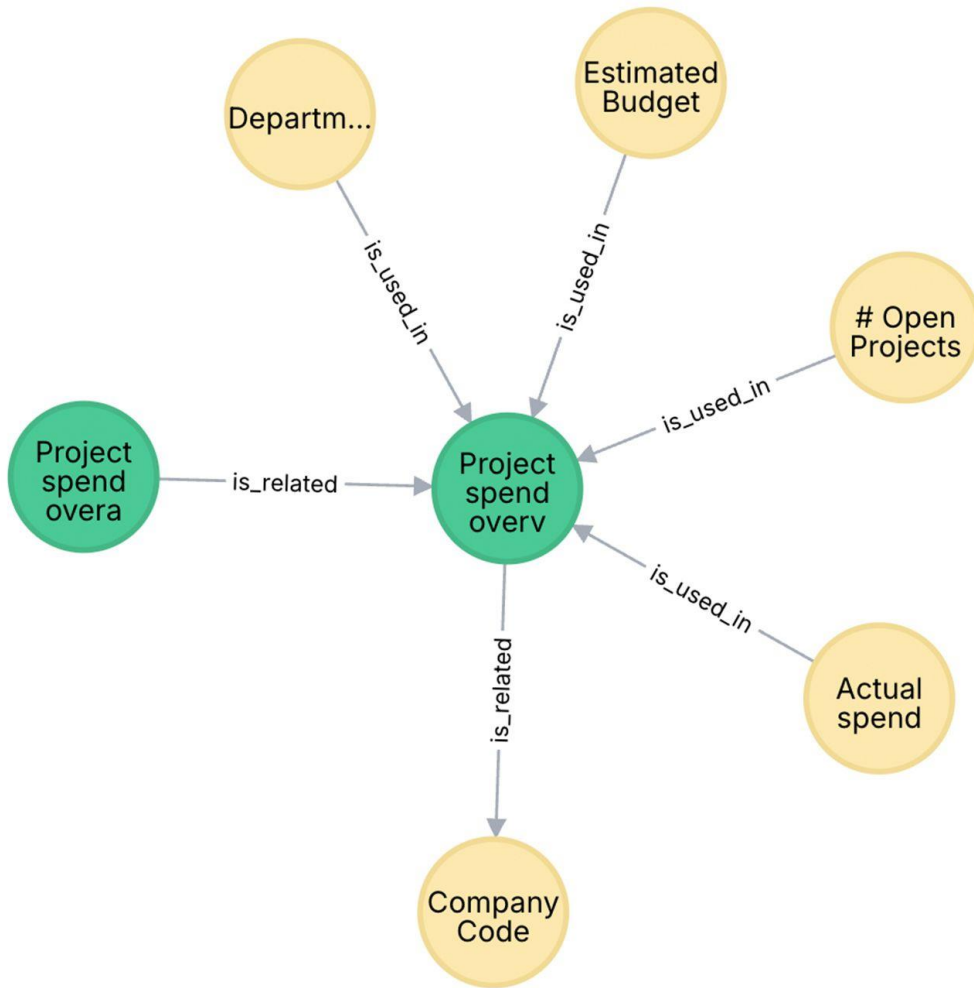
KPIs

List all Key Performance Indicators available in our organisation

MD

Marketing domain

List all Items associated to Marketing domain

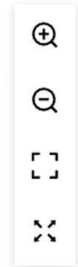


Asset Types

- Report
- Definition
- Webpage

Relation Type

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



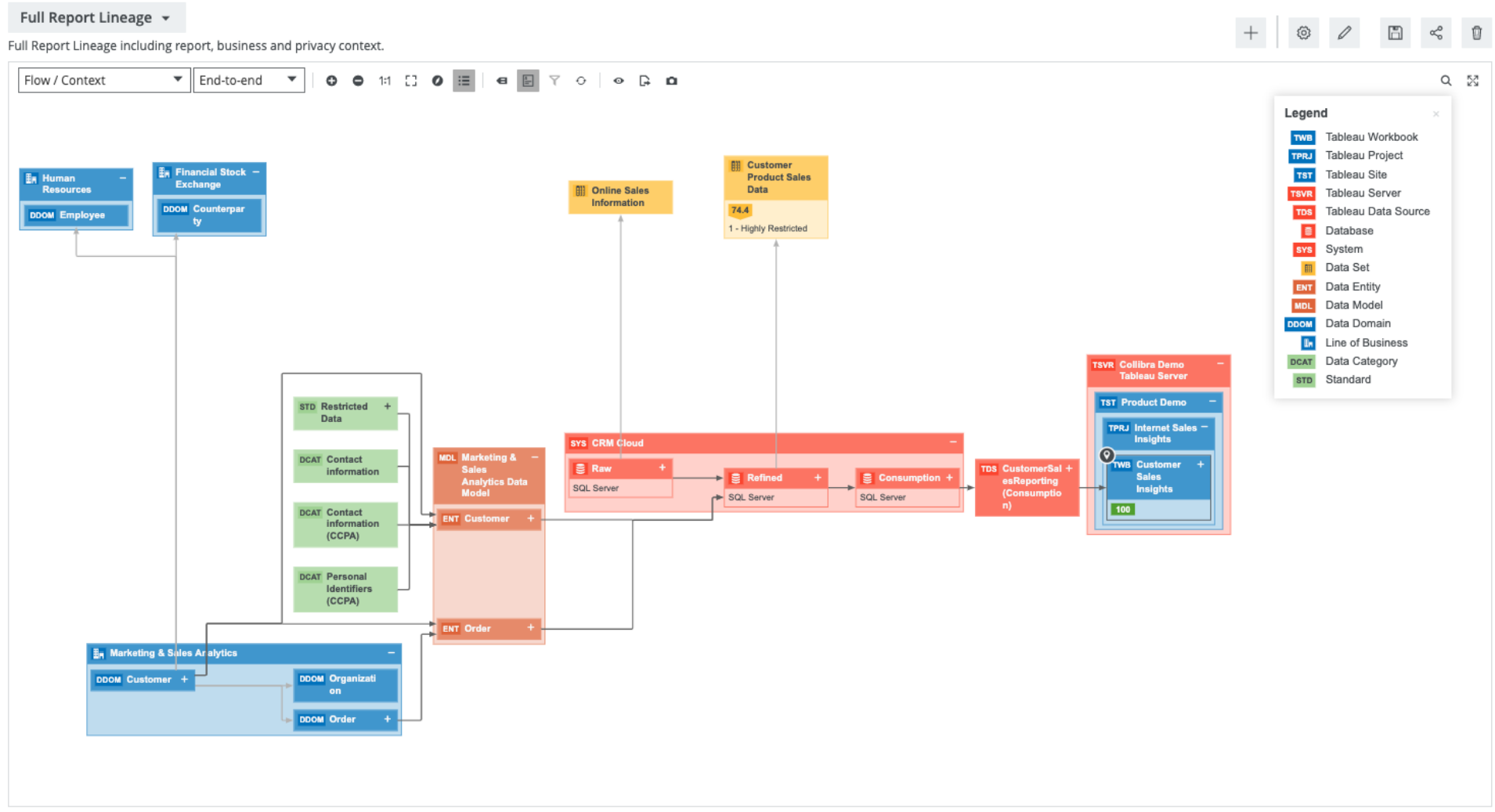


TWB Customer Sales Insights

Tableau Workbook | Accepted | ★★★★★ (1) | 0 | 1

Add to Data Basket More

- Add characteristic
- Details
- Diagram
- Pictures
- Quality
- Responsibilities
- References
- History
- Files

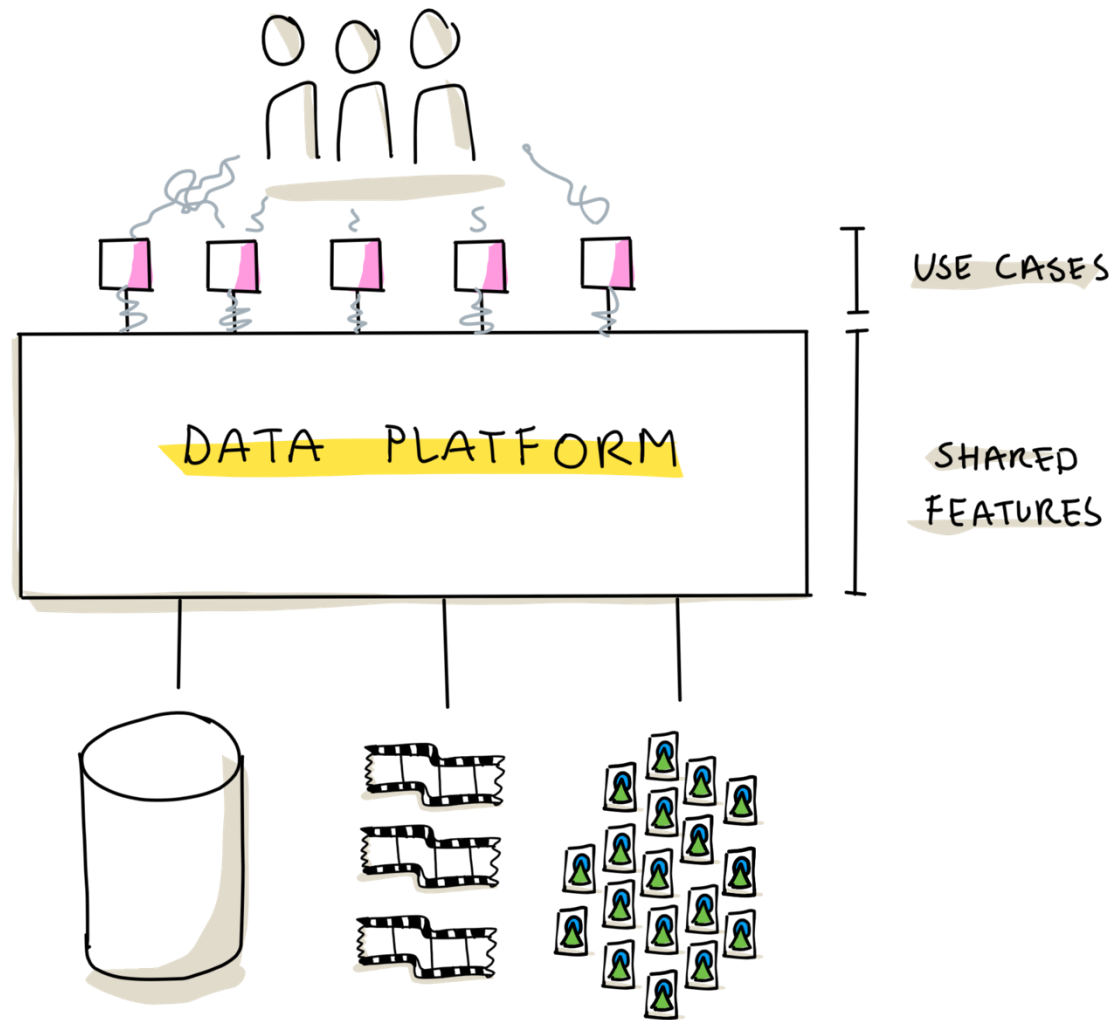
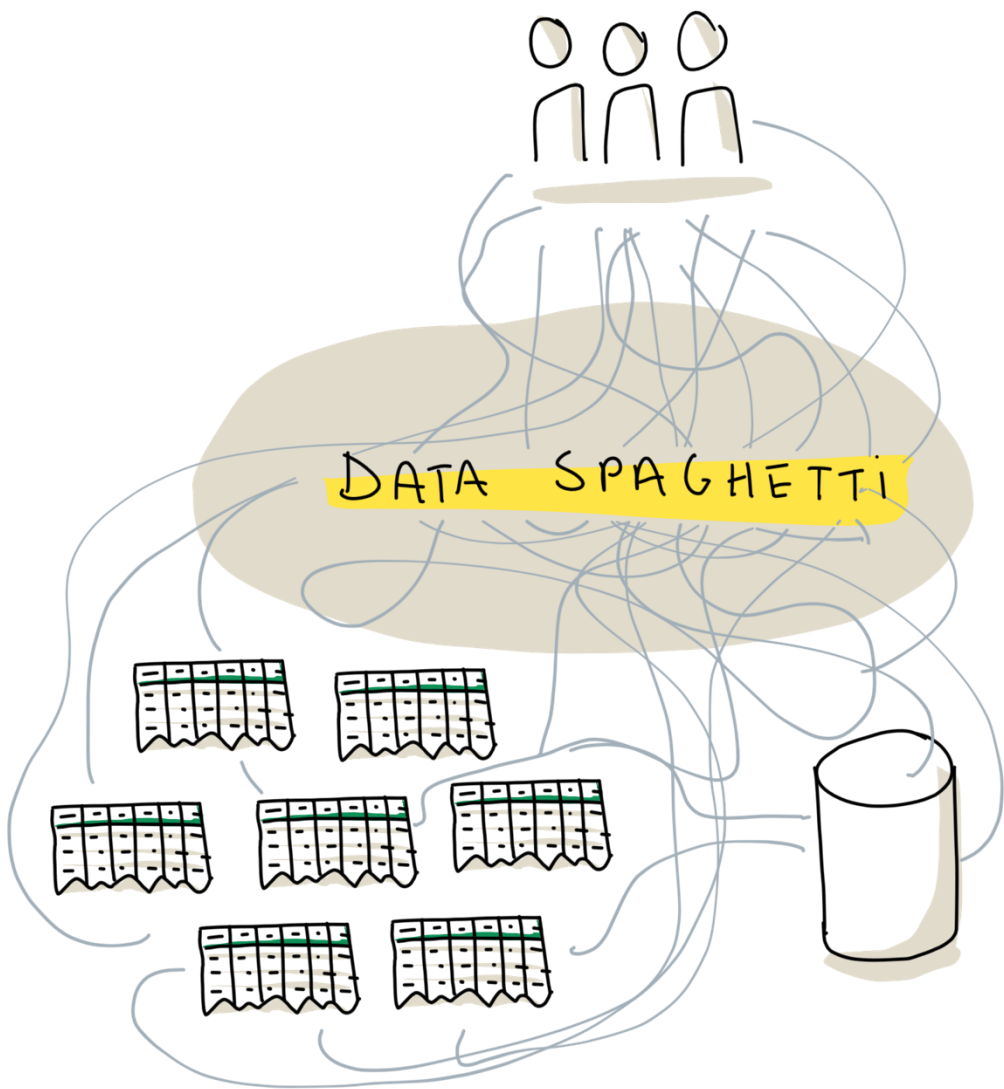






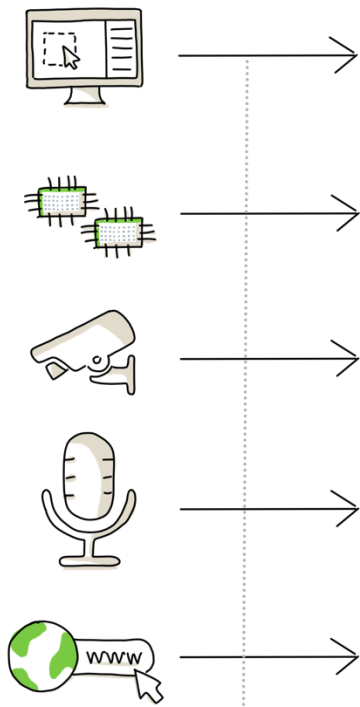
Data Management

- Defining Data
- Data Producers
- Big Data Vs
- The Bigger Picture
- Data Management
- **Data Technology**

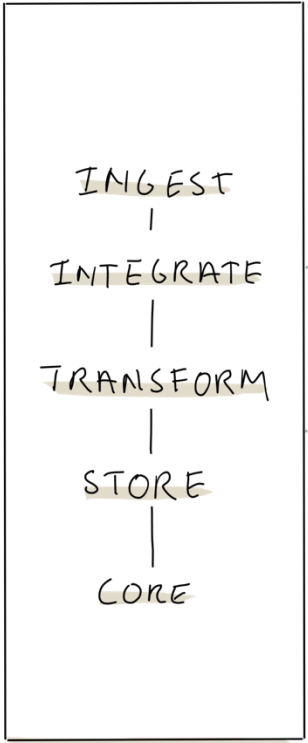




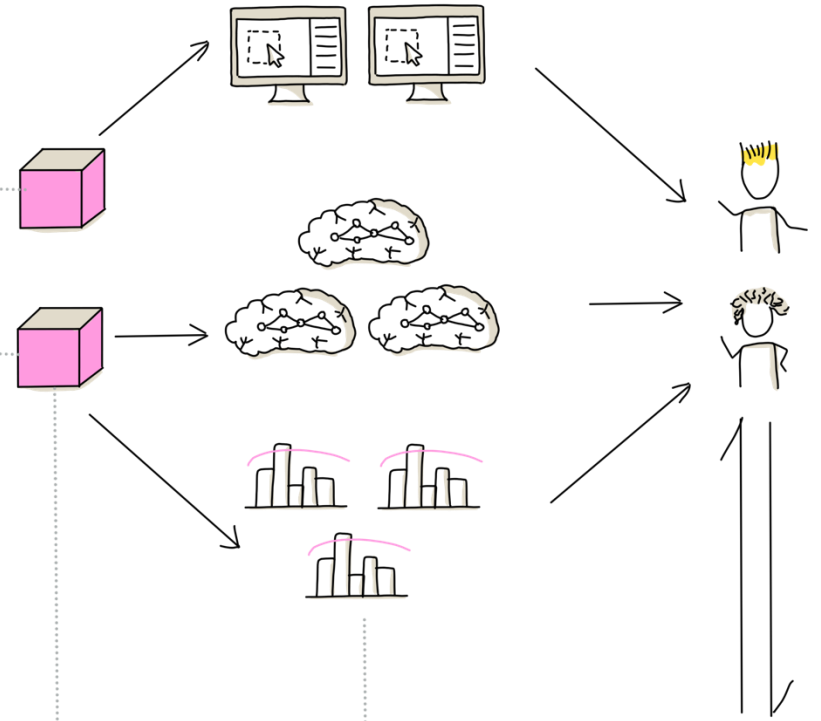
DATA PRODUCERS (SOURCES)



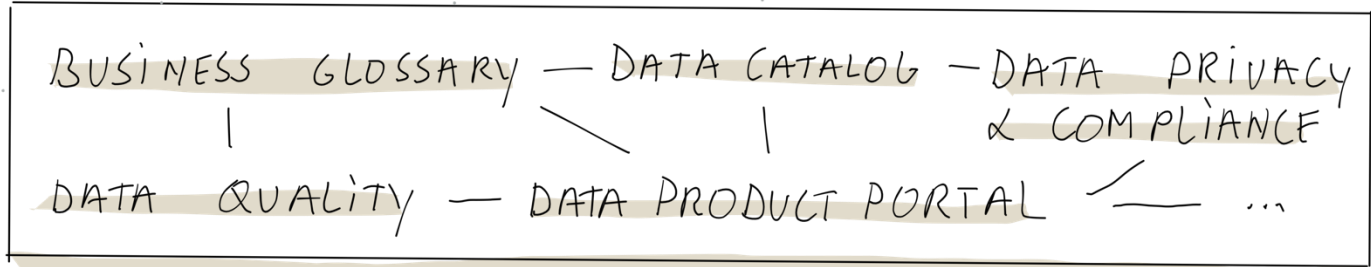
DATA PLATFORM



CONSUMPTION



VALUE



META-DATA MANAGEMENT



Modern Data Platform Architecture

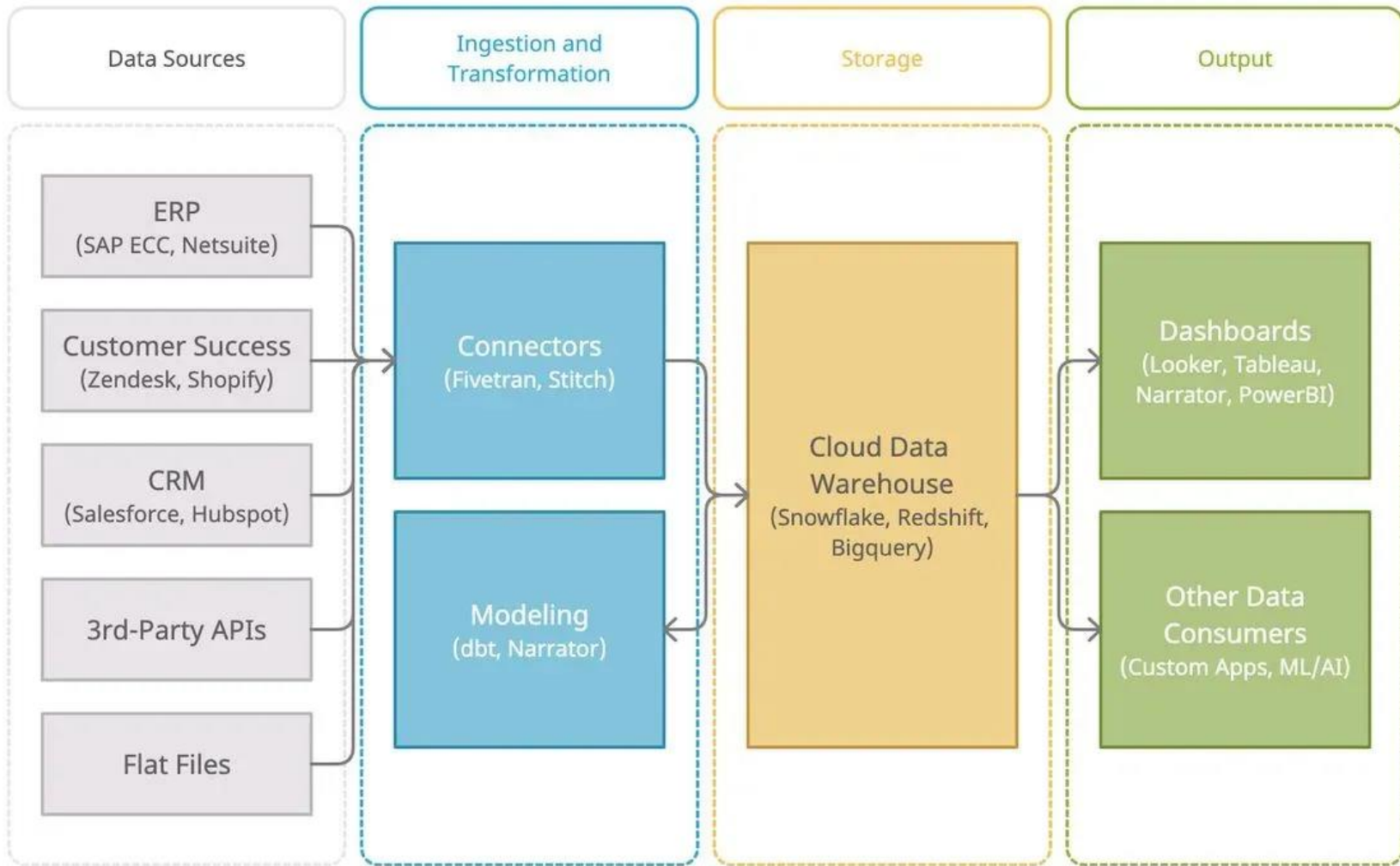
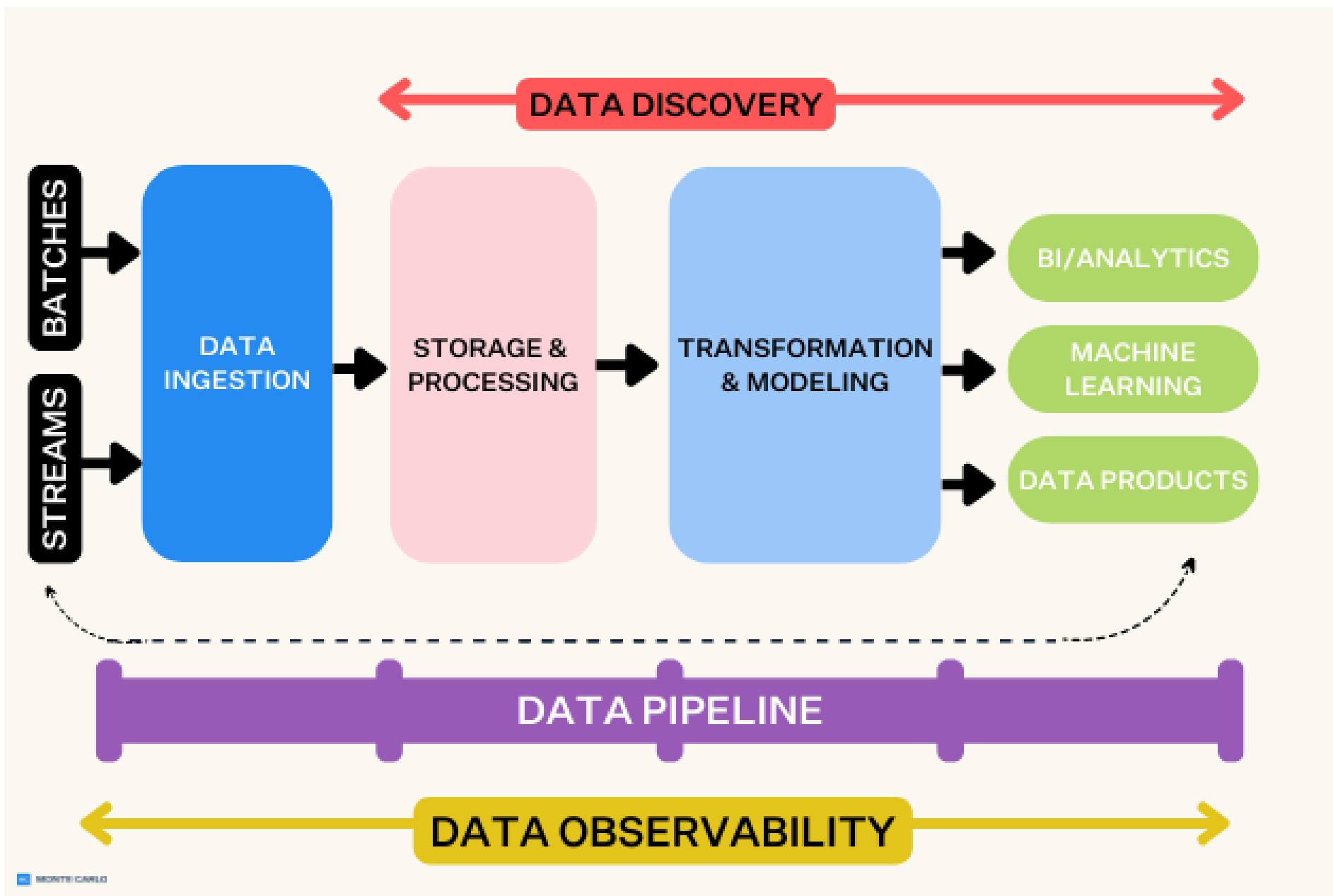
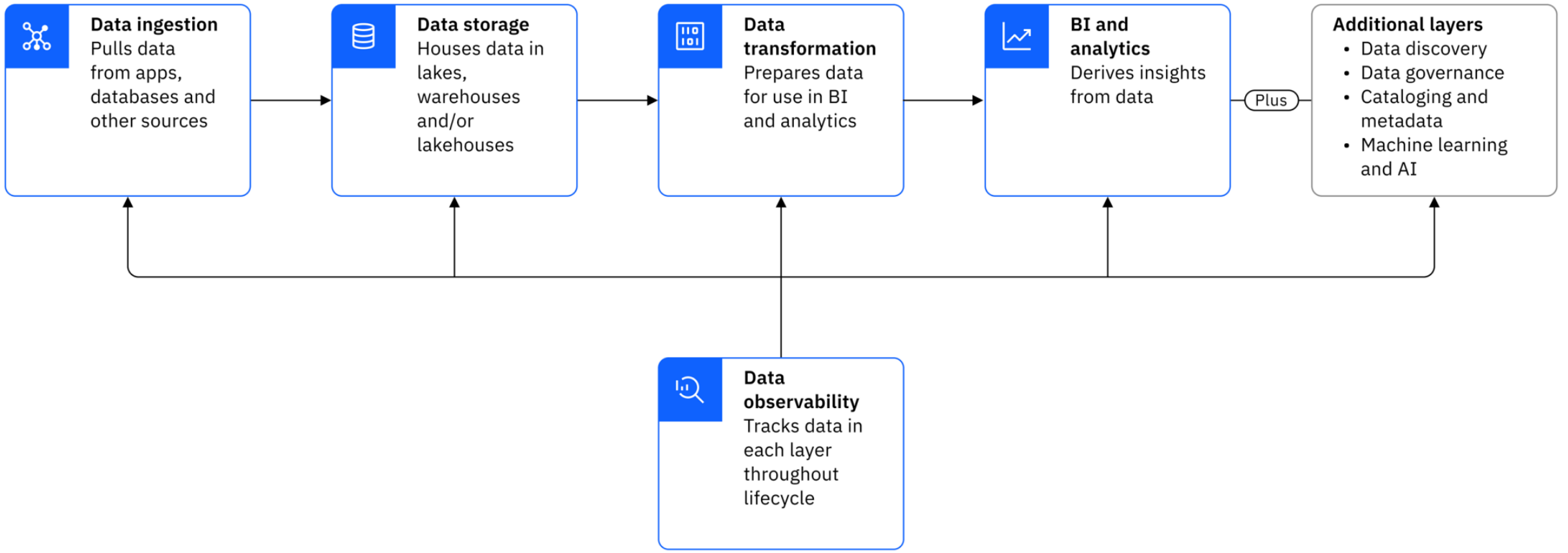
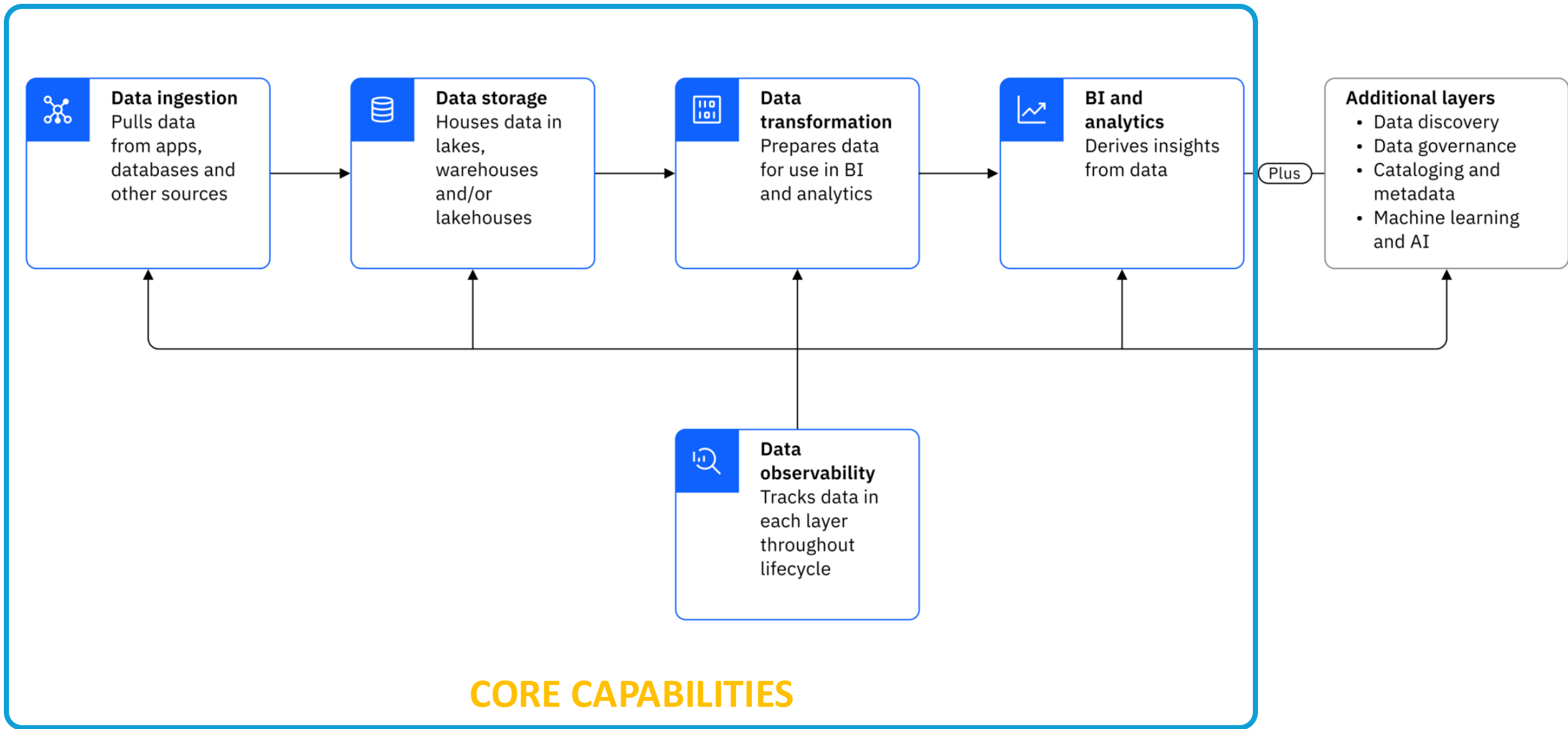
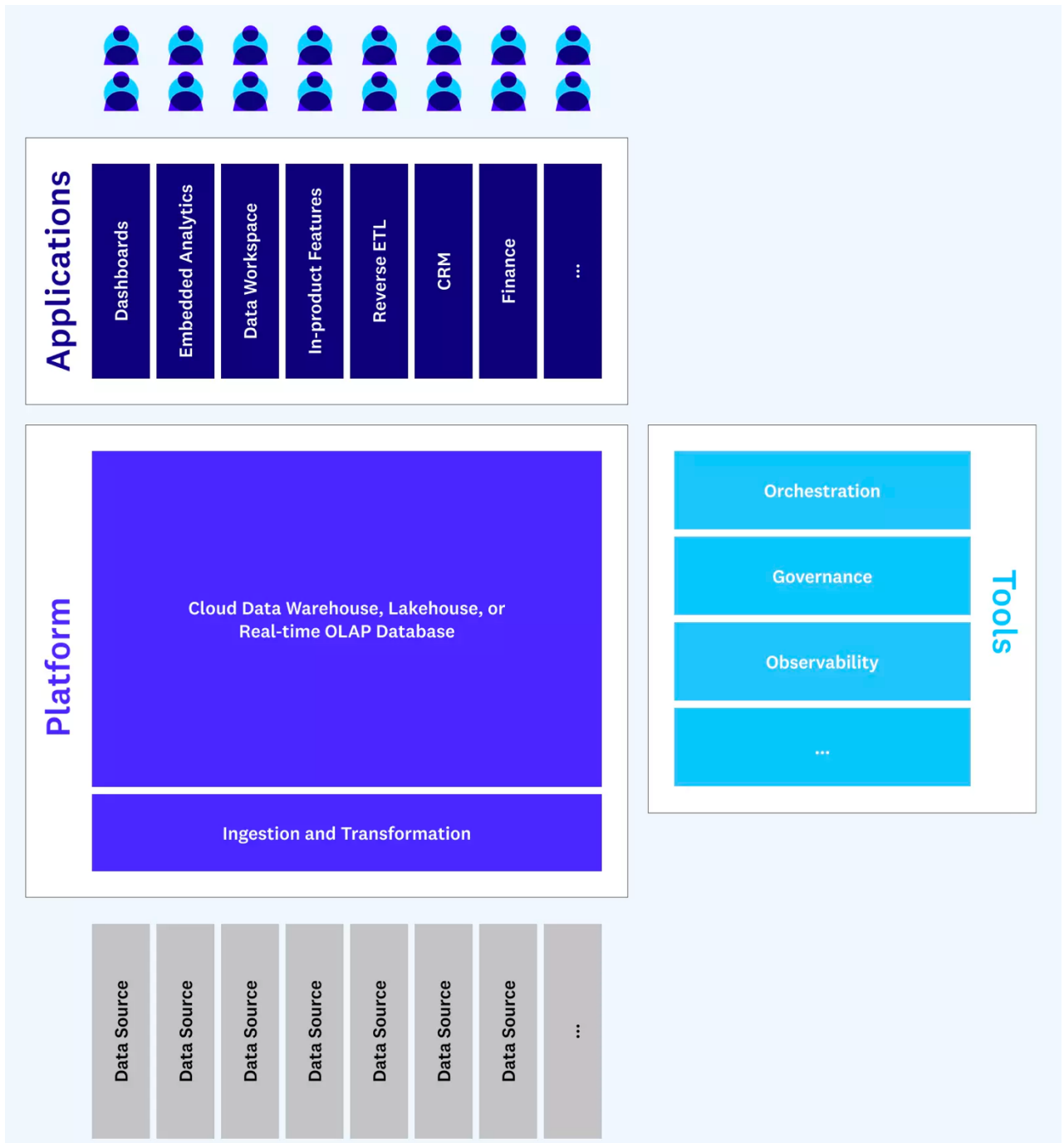


Illustration by Sanjiv Prasad









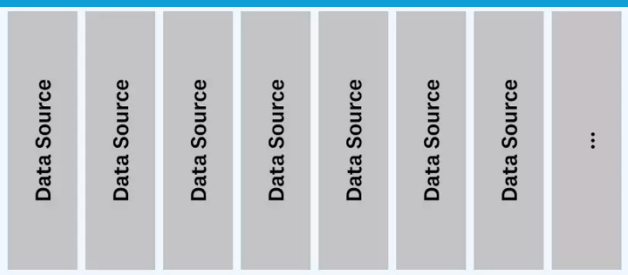
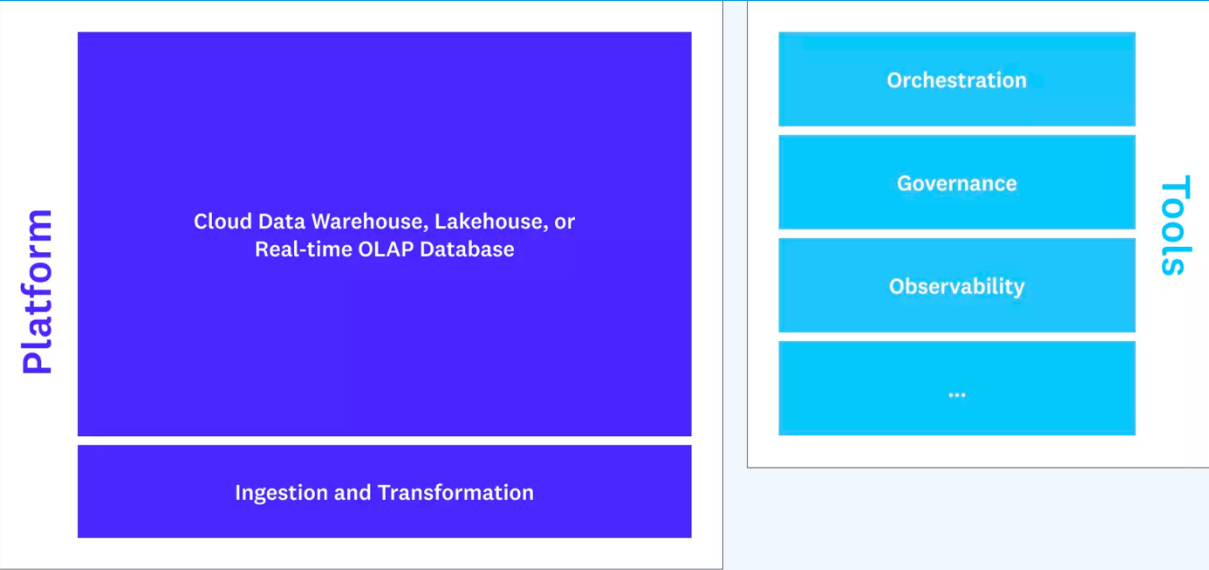
[Atlan: What is a Data Platform]



USE CASES



PLATFORM





Data Management

- Defining Data
- Data Producers
- Big Data Vs
- The Bigger Picture
- Data Management
- Data Technology