

Modern data stack











The modern data stack is one framework used to conceptualize how different data tools work together to allow a complete data journey.

Modern Data Stack

- **Infraestructure as code**

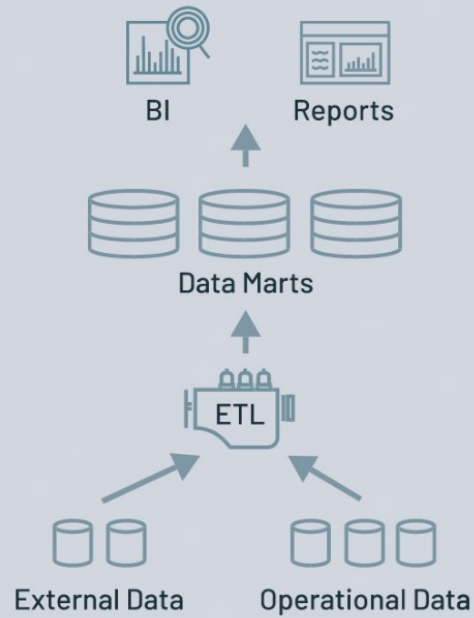


Modern Data Stack

- Infrastructure as code
- **From Datawarehouse to Lakehouse**

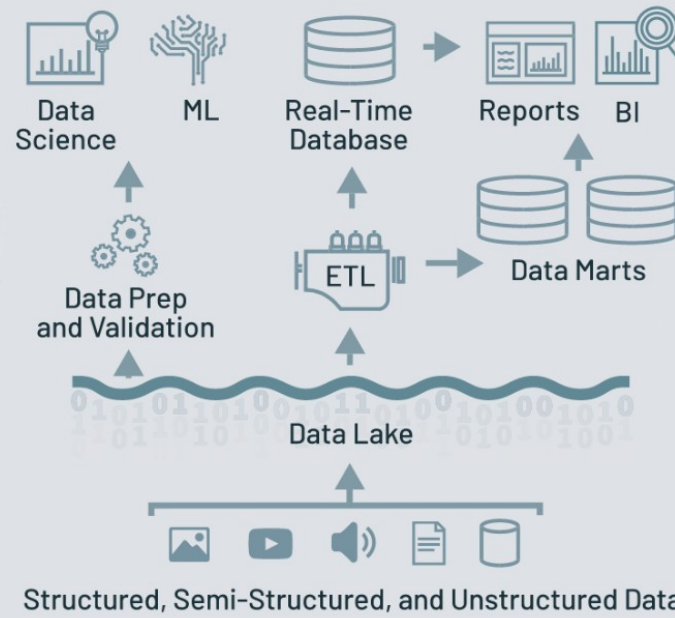
Late 1980s

DATA WAREHOUSE



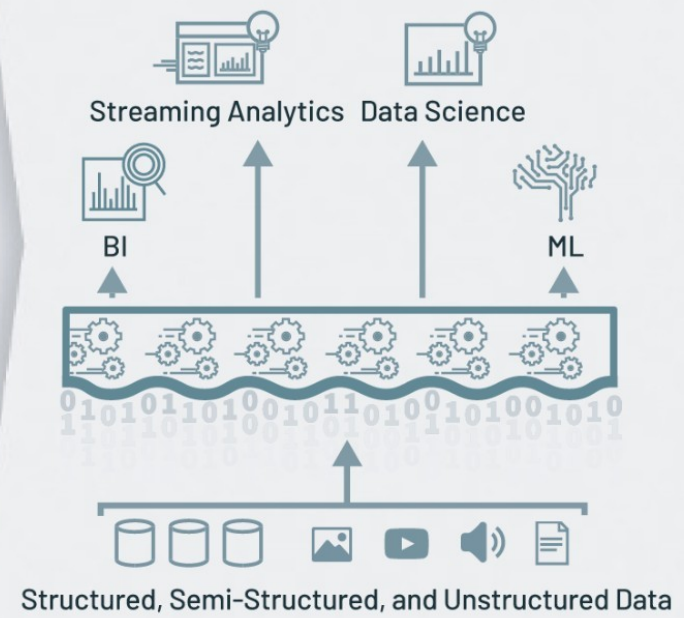
2011

DATA LAKE



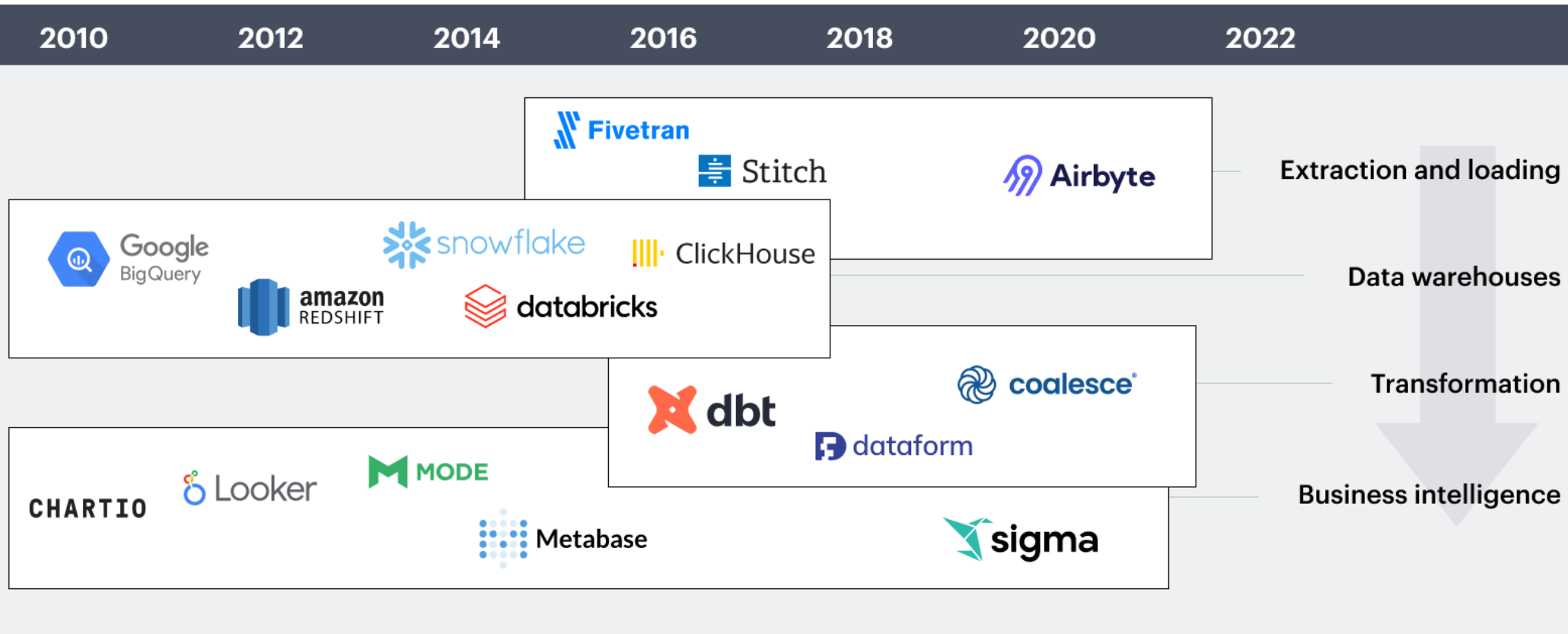
2020

LAKEHOUSE



Modern Data Stack

- Infrastructure as code
- From Datawarehouse to Lakehouse
- **From ETL to ELT**

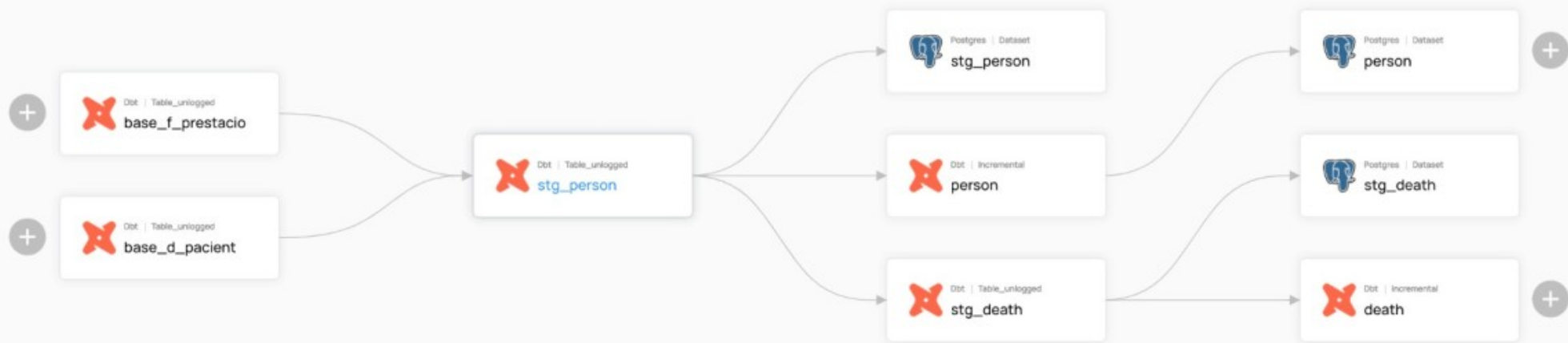


Adoption of the Modern Data Stack

Helps data teams pull data from sources (**extraction and loading**) into one place (**warehouses**), transform it into a usable form (**transformation**), and then give access with accessible tools (**business intelligence**).

Modern Data Stack

- Infrastructure as code
- From Datawarehouse to Lakehouse
- From ETL to ELT
- **Data governance and observability**



Accuracy

The degree to which the data correctly represents the entity or attribute being described.



Completeness

The percentage of missing data from a given data set.



Consistency

The absence of difference or contradiction in data irrespective of the data's source.



Validity

Invalid data affects the accuracy and completeness of a data set.



Integrity

The validity of relationships across various data entities.



Uniqueness

Ensures duplicate or overlapping data is identified and marked.

Modern Data Stack

- Infrastructure as code
- From Datawarehouse to Lakehouse
- From ETL to ELT
- Data governance and observability
- **Self-service analytics**

PRODUCTION

APPME D

Search by Title...

+

Recent

▼

≡

PERSON

Adhoc

Person analysis

Domain Visit

DRUGS (QA)

EDA: DRUG_EXPOSURE (H12O)

☆ Favorites

↓U@

▼

No items in this section.

≡ Lists

+

▼

OMOP

>

📄 All My Docs

↓U@

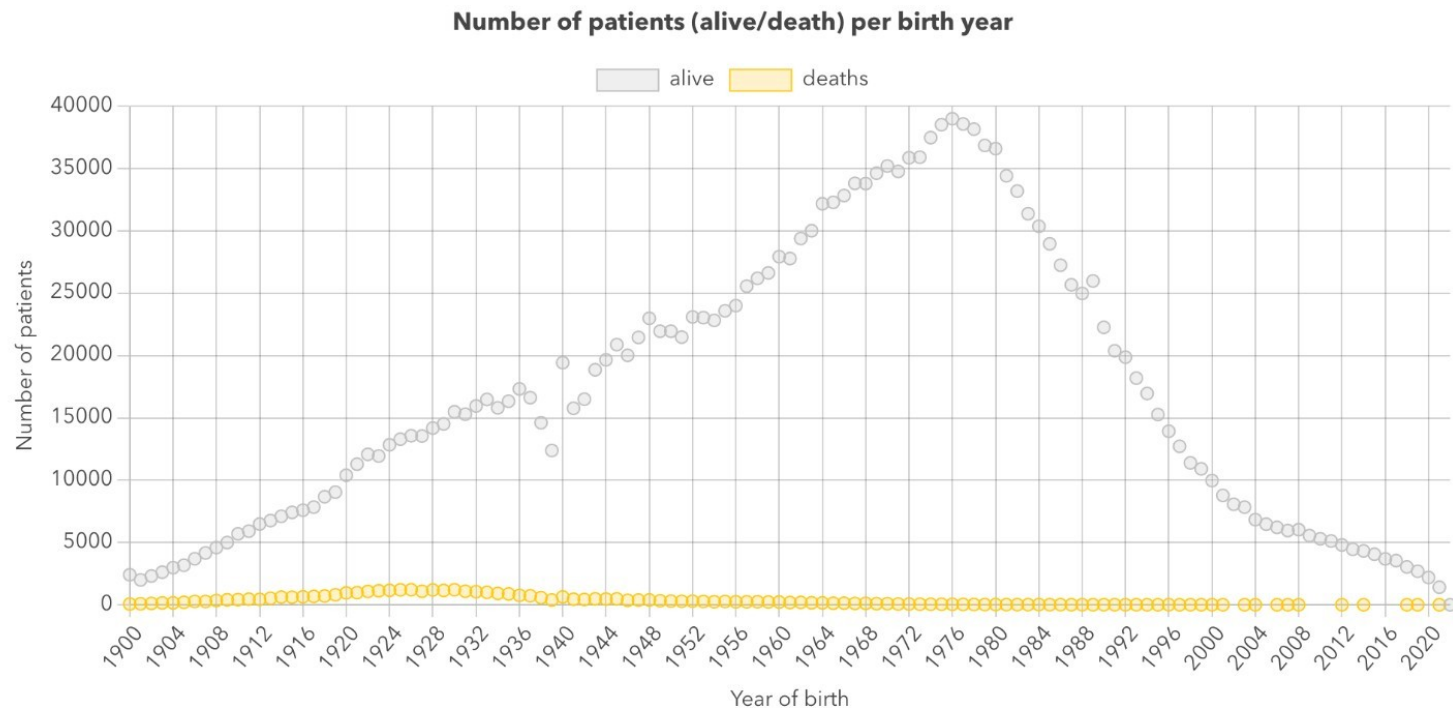
▼

Untitled

Search

EXECUTION 90 ● ▼
STATEMENT 1 OUT OF 1 ▼

CONFIG CHART ⓘ



Governance



DataHub

Sources

Applications & APIs



Databases



Files



Events



Ingestion



Storage



PostgreSQL



DuckDB



ClickHouse



dremio

Data analytics



Querybook



Apache
Superset™



Data quality



great
expectations



elementary

Orchestration



Apache

Airflow



dagster

Transformation



dbt



APACHE
Spark

Practice time?



bsc-health-data / pycones-23-modern-data-stack

Q Type ↗ to search

<> Code

Pull requests

Actions

Projects

Wiki

Security

Insights

Settings

pycones-23-modern-data-stack Public

Edit Pins ▾

Watch 0

forked from [bsc-health-data/pydatalondon23-modern-data-stack](#)

main ▾

1 branch

0 tags

Go to file

Add file ▾

<> Code ▾

This branch is [32 commits ahead](#) of bsc-health-data:main.

Contribute ▾

Sync fork ▾

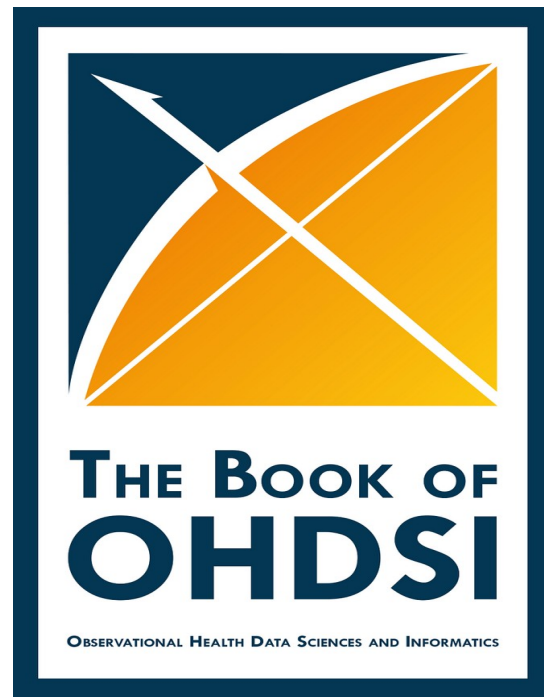
alabarga Update README.md 5b0283a yesterday 108 commits

database	Add database DOckerfile	last year
datahub	Add datahub	last year
dremio	Add files	last year
meltano	Add DBT project	5 months ago
obevo	Add obevo schemas	last year

<https://github.com/bsc-health-data/pycones-23-modern-data-stack/>

Recursos

- Sitio web: <https://github.com/alabarga/omop-cdm-course>
- Diapositivas
- El Libro de OHDSI
- EHDEN academy
- OHDSI workshops
- Atlas tutorials
- Entorno local (docker)



You are free to:

Share — copy and redistribute the material in any medium or format

Under the following terms:



Attribution - You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.



NonCommercial - You may not use the material for commercial purposes. |



NoDerivatives - If you remix, transform, or build upon the material, you may not distribute the modified material.