# Data Governance

Kristo Raun
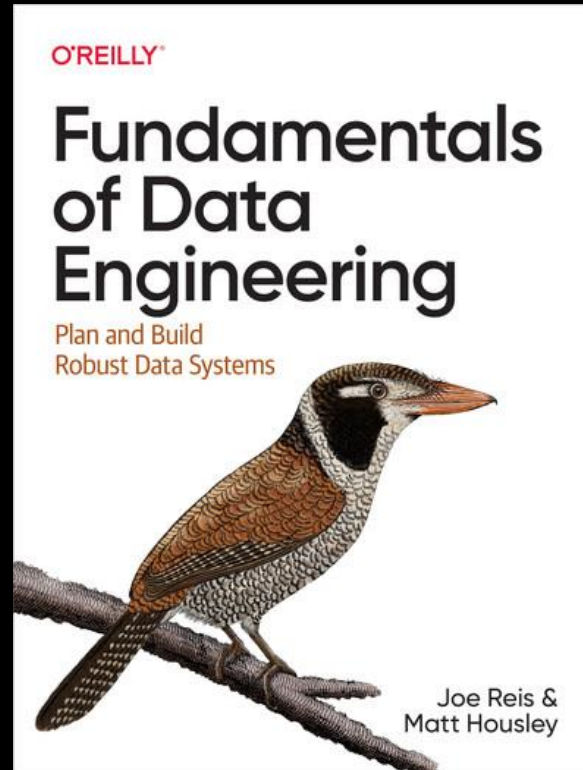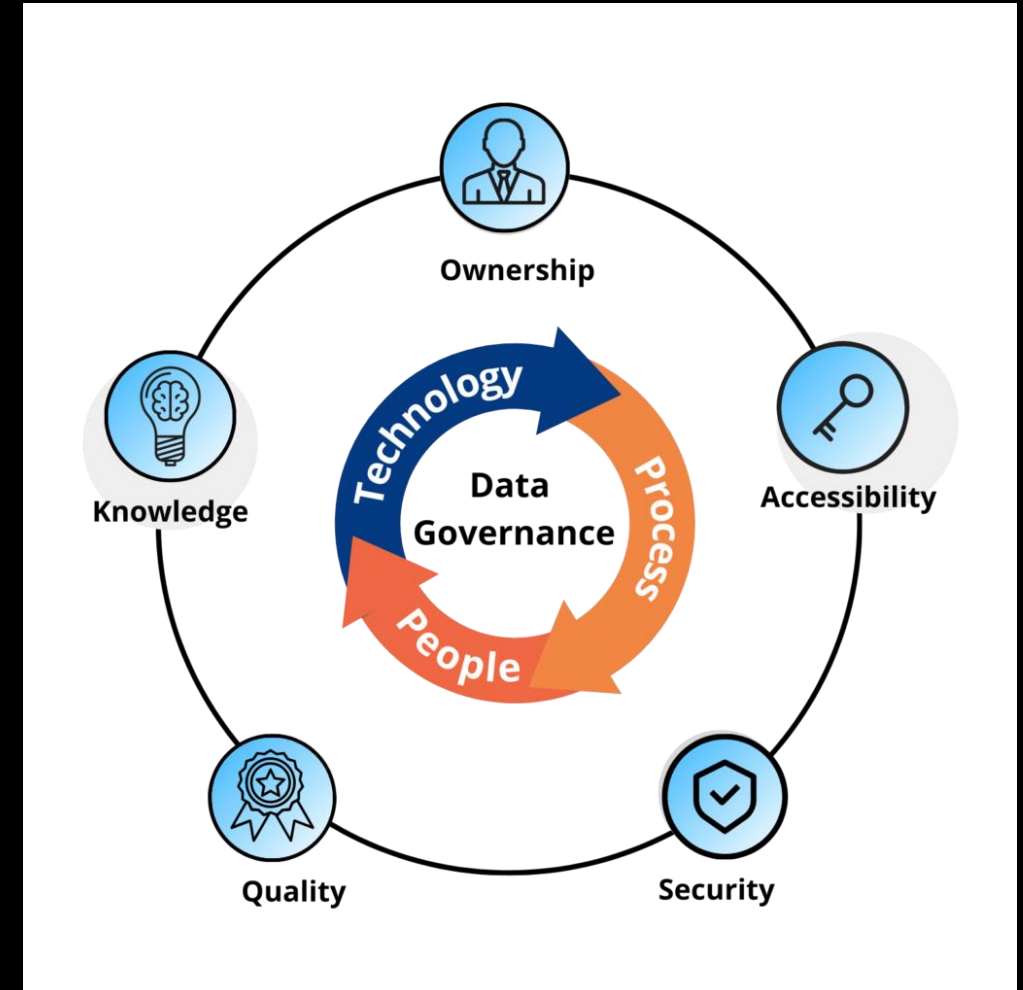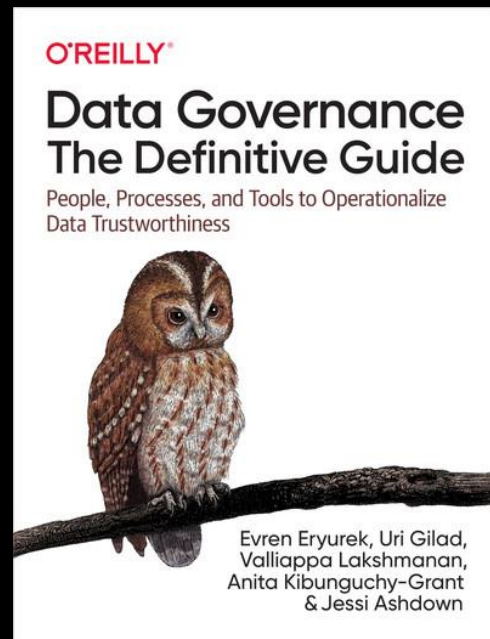
Data Engineering 2024 Fall

# Reading

- Chapter II
  - Data Governance

# Data Governance

"Data governance is, first and foremost, a data management function to ensure the **quality**, **integrity**, **security**, and **usability** of the data collected by an organization."



O'REILLY®

**Data Governance**
**The Definitive Guide**

People, Processes, and Tools to Operationalize
Data Trustworthiness

Evren Eryurek, Uri Gilad,
Valliappa Lakshmanan,
Anita Kibunguchy-Grant
& Jessi Ashdown

# Agenda

- Data Quality
- Data Usability (Data Discoverability)
- Quiz session

# Data Quality



1. Accuracy (Õigsus), e.g. Does the data entry match real-world values?
   - Correct: Name: "Tõnu," Gender: "M" (matches actual values).
   - Incorrect: Name: "T6nu," Gender: "N" (typo, incorrect gender).

2. Completeness (Täielikkus), e.g. Are all required fields filled?
   - Correct: Customer record with all fields (Name, Email, Phone) populated.
   - Incorrect: Customer record missing "Email" field.

Andmekvaliteedi juhis, mai 2023
https://www.kratid.ee/en/juhised

3. Timeliness (Ajakohasus), e.g. Is the data available when needed?
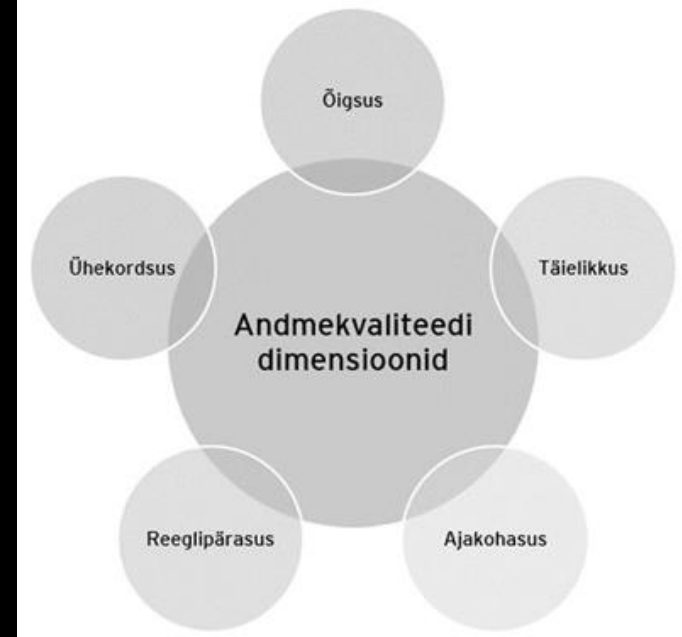   - Correct: Delivery status updated in real time for tracking.
   - Incorrect: Delivery status updated only weekly, causing delays in tracking.

4. Orderliness (Reeglipärasus), e.g. Does the data follow the agreed format?
   - Correct: Date of birth recorded as "2024-11-09" (ISO format).
   - Incorrect: Date of birth recorded as "11/9/24" (ambiguous, non-standard).

5. Uniqueness (Ühekordsus), e.g. Are there duplicate entries in the dataset?
   - Correct: Only one record exists per customer (e.g., John Doe has one entry).
   - Incorrect: Same customer has multiple records (e.g., "John Doe" and "J. Doe" for the same individual).
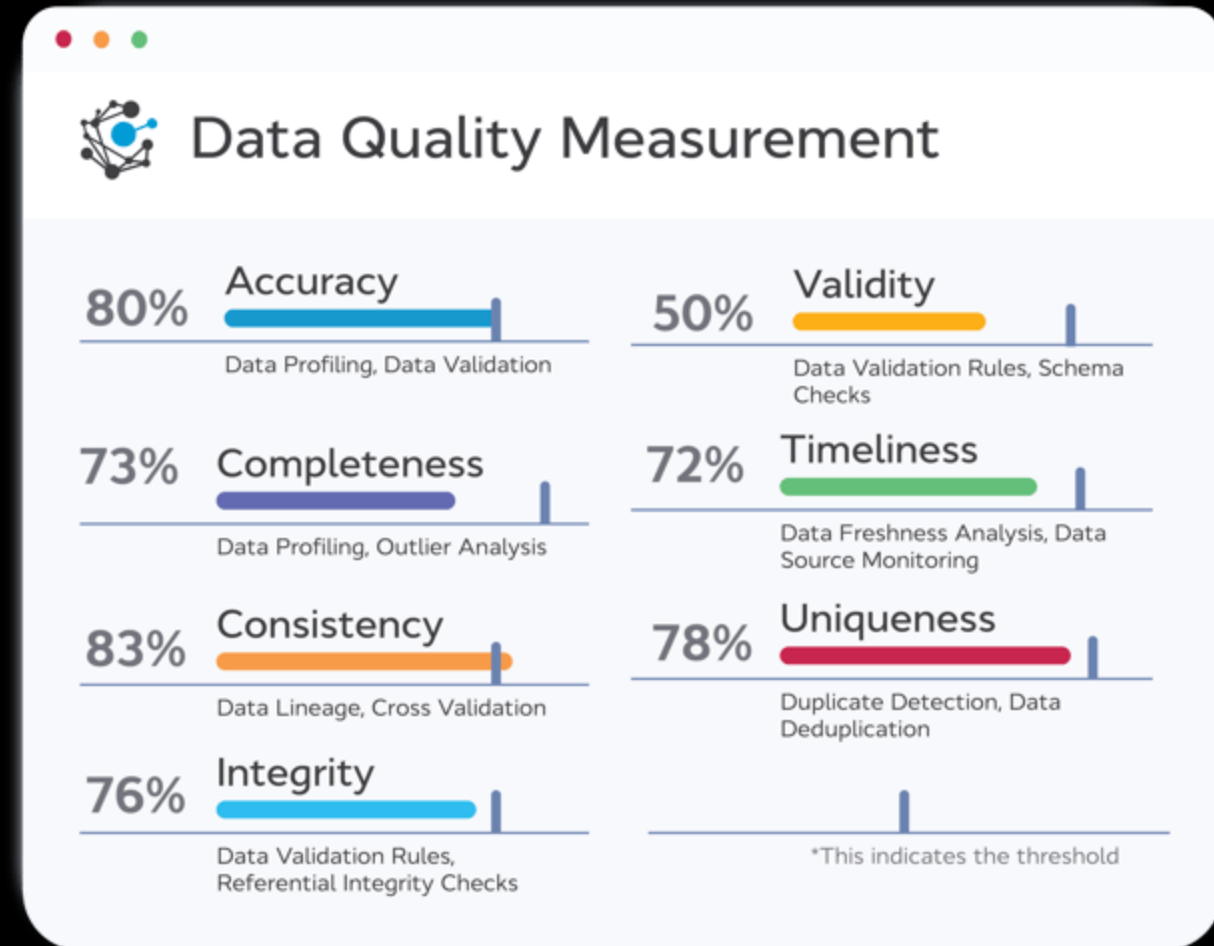
# Data Quality Metrics

- Perfect Data Quality is impossible
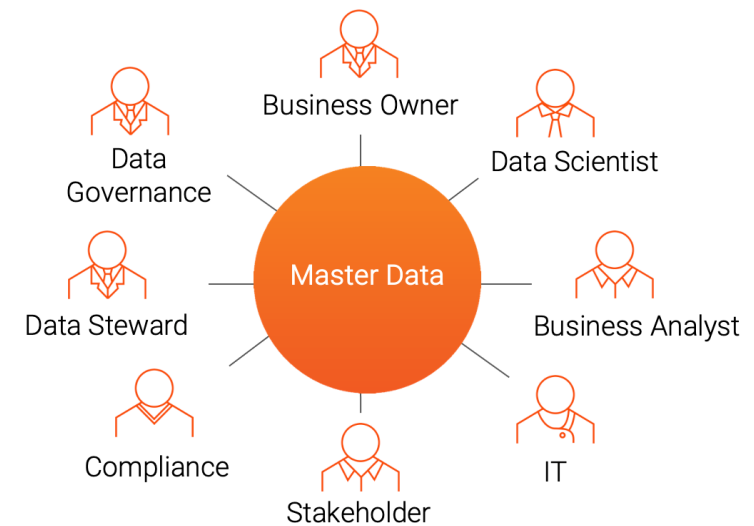  - How do we evaluate data quality?
  - What level is acceptable?



https://www.ovaledge.com/blog/data-quality-metrics

# Data Quality: Master Data Management

- Consistent entity definitions
  - Standard formats
  - Single source of truth
    (for key entities)
- E.g.:
  - Customer
  - Product
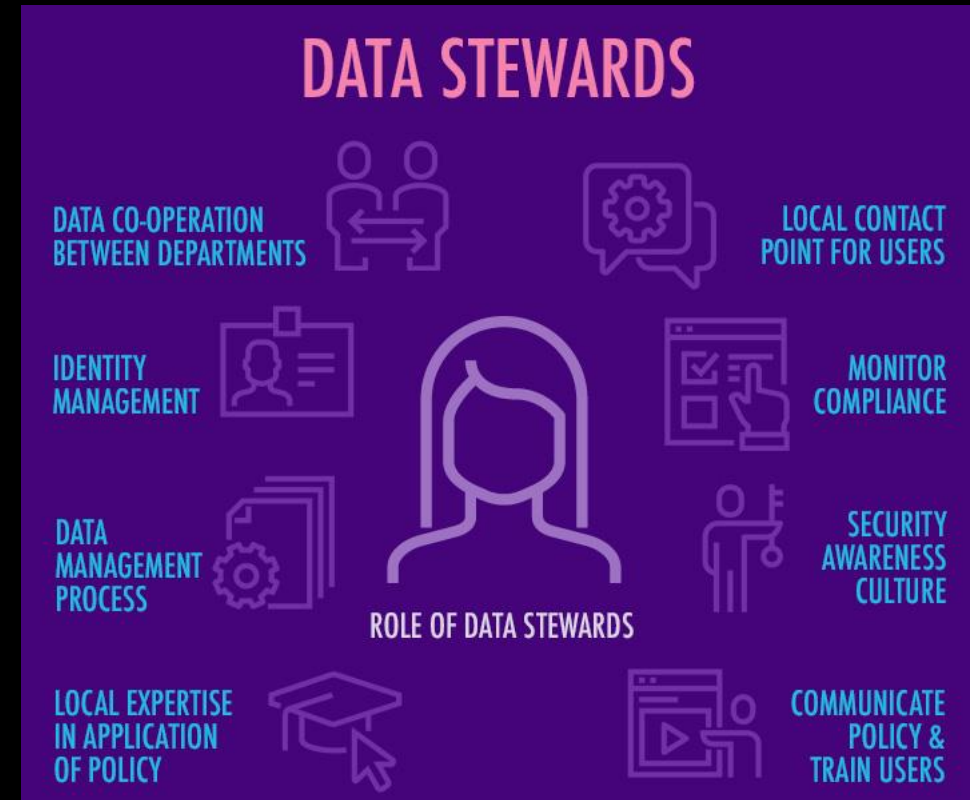  - Location



Master Data Management Explained

Master data management involves creating a single master record for each person, place, or thing in a business, from across internal and external data sources and applications.

https://www.informatica.com/resources/articles/what-is-master-data-management

# Data Quality: Data Stewardship

- Accountability: who is "in charge" of the data?

- Data Steward: Managing, overseeing, and enforcing data governance policies within an organization.
    - Data Quality Monitoring
    - Data Compliance
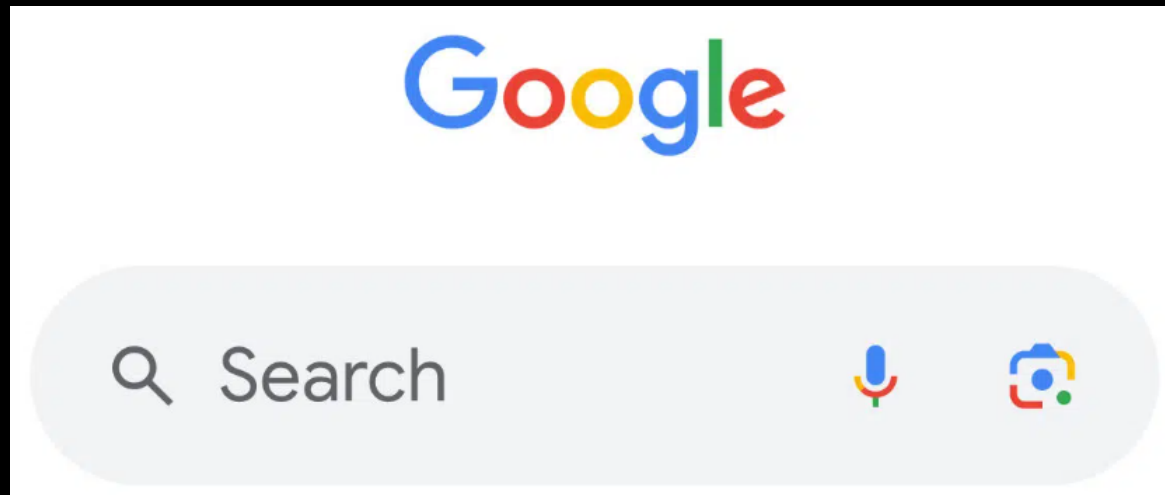    - Helping data users



DATA STEWARDS

DATA CO-OPERATION BETWEEN DEPARTMENTS

LOCAL CONTACT POINT FOR USERS

IDENTITY MANAGEMENT

MONITOR COMPLIANCE

DATA MANAGEMENT PROCESS

SECURITY AWARENESS CULTURE

ROLE OF DATA STEWARDS

LOCAL EXPERTISE IN APPLICATION OF POLICY

COMMUNICATE POLICY & TRAIN USERS

# Agenda

- ~~Data Quality~~

- Data Usability (Data Discoverability)

- Quiz session

# Data Usability & Discoverability

- Data must be available (usable) and discoverable.
  - Data serves a purpose: "End users should have quick and reliable access to the data they need to do their jobs."
  - "They should know where the data comes from, how it relates to other data, and what the data means."

# Metadata

- What is metadata?
  - Data about data
  - Why do we need metadata?
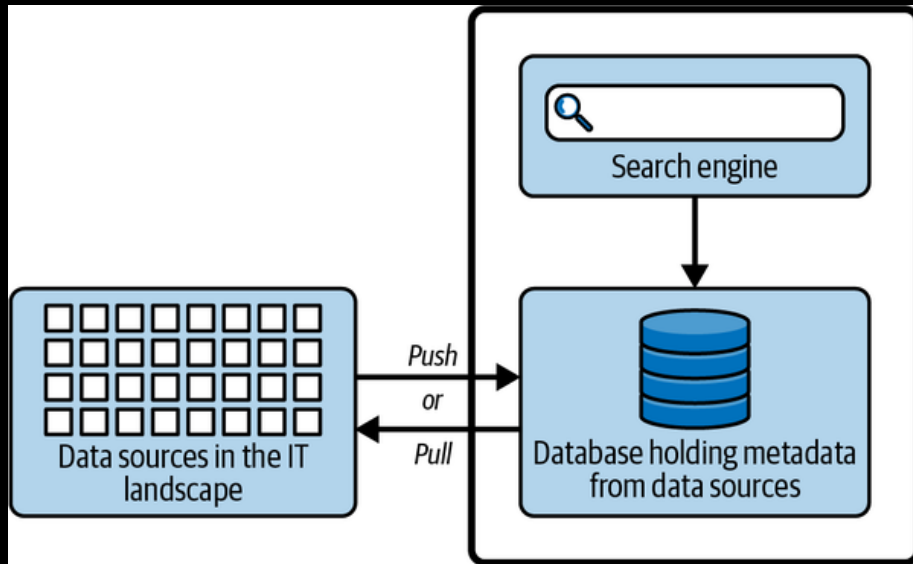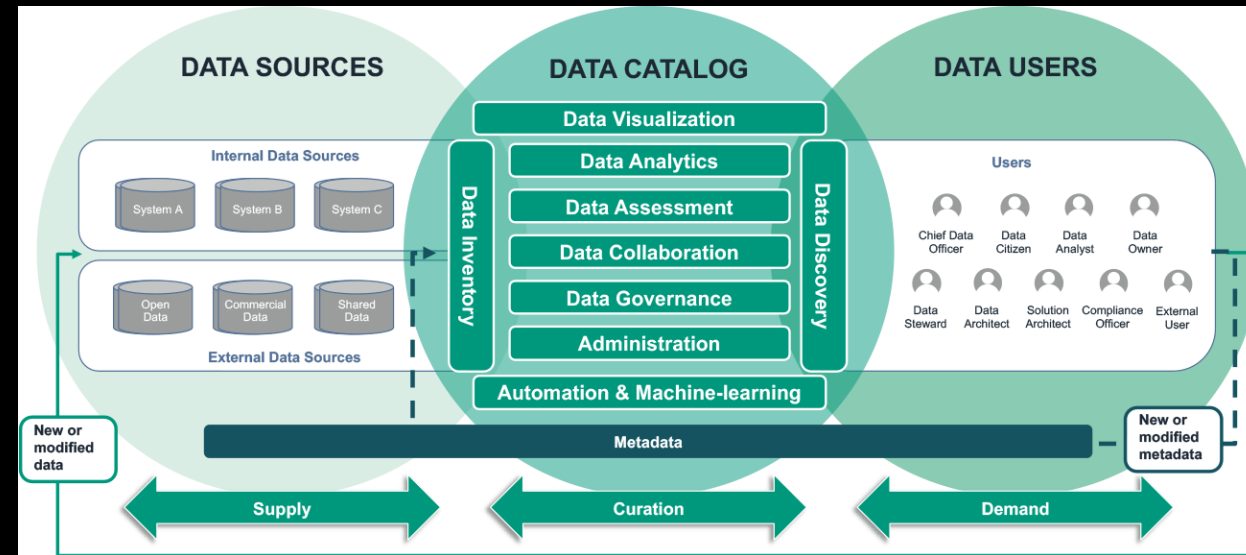
organized



unorganized

# Metadata categories

| Category | Description | Examples |
|---|---|---|
| Business metadata | Relates to how data is used within the business context, including definitions, rules, and ownership. | • Data Definitions (who is a "customer"?)<br>• Data Rules (phone number must be 8 digits)<br>• Data Ownership (registered by, owned by, …) |
| Technical metadata | Describes the technical aspects of data, including structure, lineage, and system interactions. | • Schema Metadata (Database schemas, table structures)<br>• Data Lineage (Tracking data origin and transformations)<br>• Pipeline Metadata (Workflow schedules, system dependencies) |
| Operational metadata | Captures information about the operations and processes that handle the data, including performance and error logs. | • Process Statistics (Job IDs, execution times)<br>• Runtime Logs (Application logs, error messages)<br>• Process Status (Success, fail) |
| Reference metadata | Used to classify and standardize other data, often serving as lookup data for consistency. | • Geographic Codes (EE, LV, LT)<br>• Units of Measurement (g, kg, tk)<br>• Internal Codes (ATI) |

# Data Catalog

- An organized inventory of the data in your company.
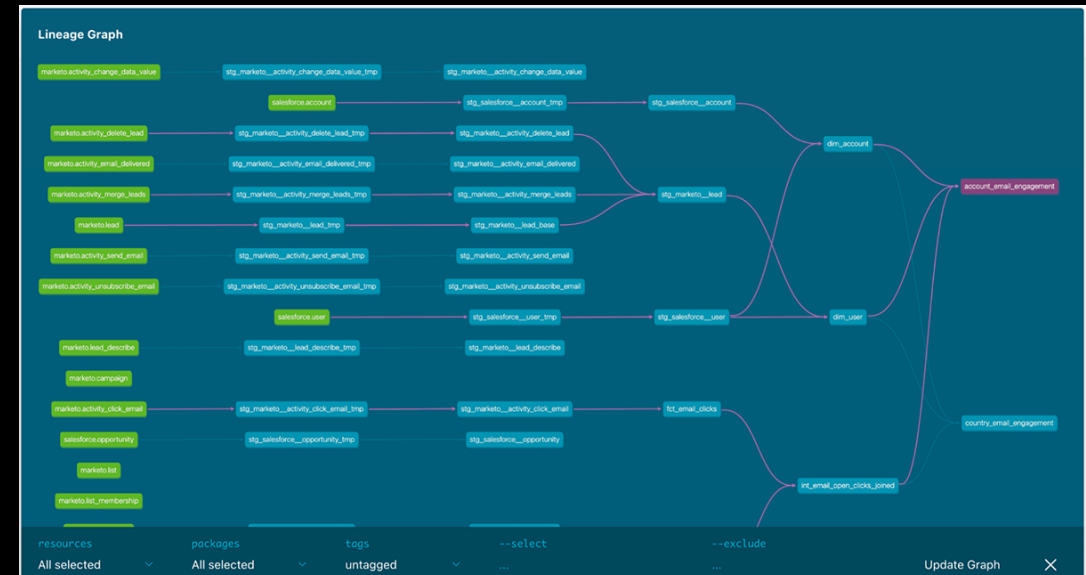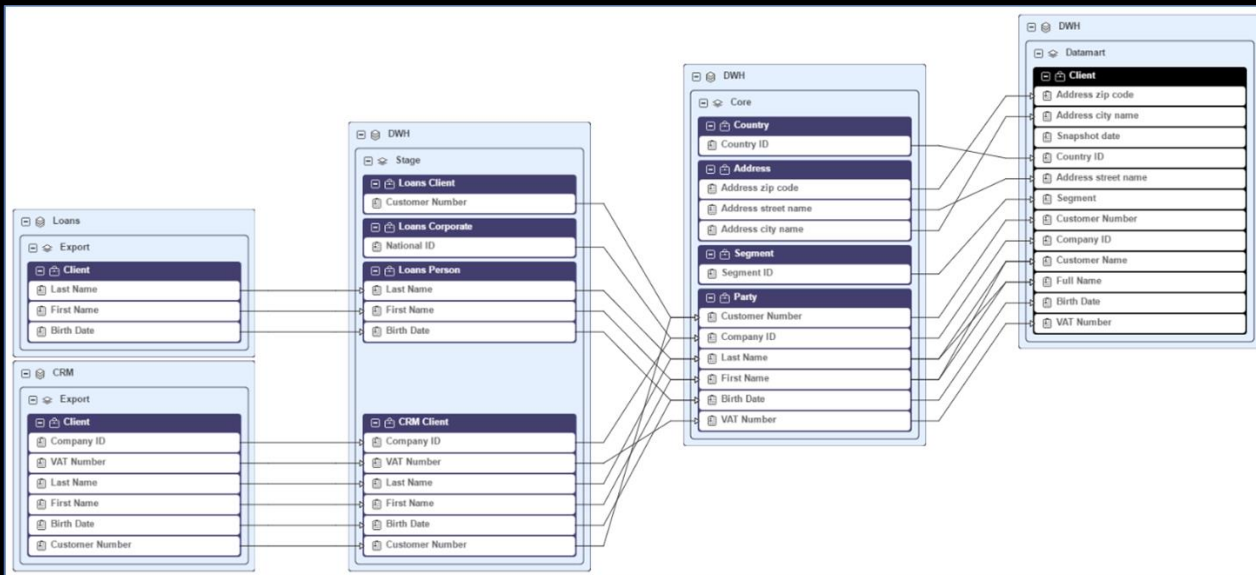


The Enterprise Data Catalog, Ole Olesen-Bagneux

hevodata.com

# Data Lineage

- Audit trail of data through its lifecycle

- Tracking across systems
  - Origins
  - Transformations
  - Consumers of data

# Agenda

- ~~Data Quality~~

- ~~Data Usability (Data Discoverability)~~

- Quiz session

# Further reading

- Chapter II
    - Data Governance
- Book
    - The Enterprise Data Catalog
- Article
    - https://www.datacamp.com/cheat-sheet/data-governance-fundamentals-cheatsheet