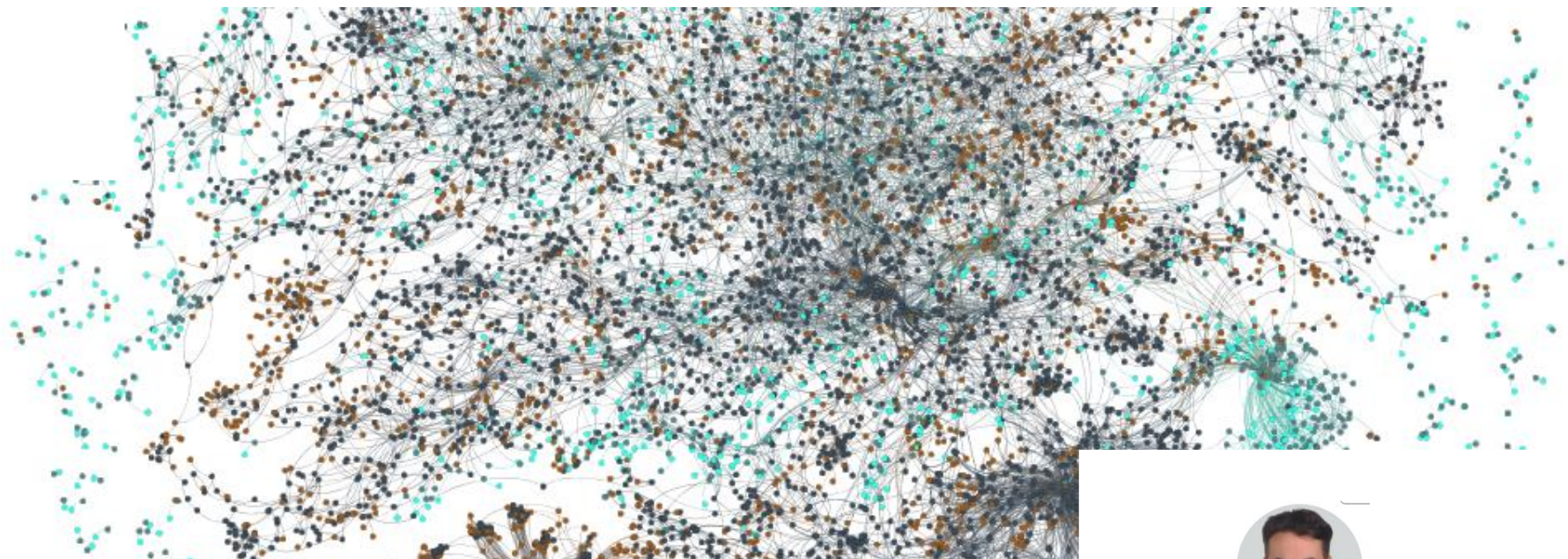# ETL, Data Management and Data Lineage

a.k.a. when life gives metadata, make graph analysis out of it



**AAI Swiss meetup**

**accenture**

Luca Perrozzi

10.11.2020

DISCLAIMER: the humor contained in this talk is based on real life experience but is not meant to provide a general representation of the professional figures involved! ☺

# The quest for data warehouses and the problems STILL to be solved

- Large financial institutions, as every other organization nowadays, produce and own **large amount of data that need to be used to take better and data-driven business decisions**

- **The technical problems implied are typically complex and challenging to solve**:
  - Fast growing data O(PB)
  - Legacy and on-premise architectures (COBOL is still a thing!)
  - Large variety of etherogeneous systems to be integrated
  - Strict security levels must be enforced

# The quest for data warehouses
# and the problems STILL to be solved

- In the last years we witnessed an **explosion of tools and architectures to address these technical problems**, and solutions are becoming more and more standardized, robust and scalable
  - Just think of Hadoop, plethora of OLTP/OLAP DBs, Cloud vendors, …

- The **real problems** (always) arise when the **technical solutions and operations need to co-exist with the business** expectations, evolving (conflicting) requirements and top-down decisions
  - Ever tried to convert business language to technical solution and viceversa?
  - Technical and business priorities are hardly overlapping
  - Communication is always an issue, especially in large organizations
  - Difficult trade-offs of **speed vs security vs costs vs maintainability**
  - «One tool to solve all problems» un-realistic promises

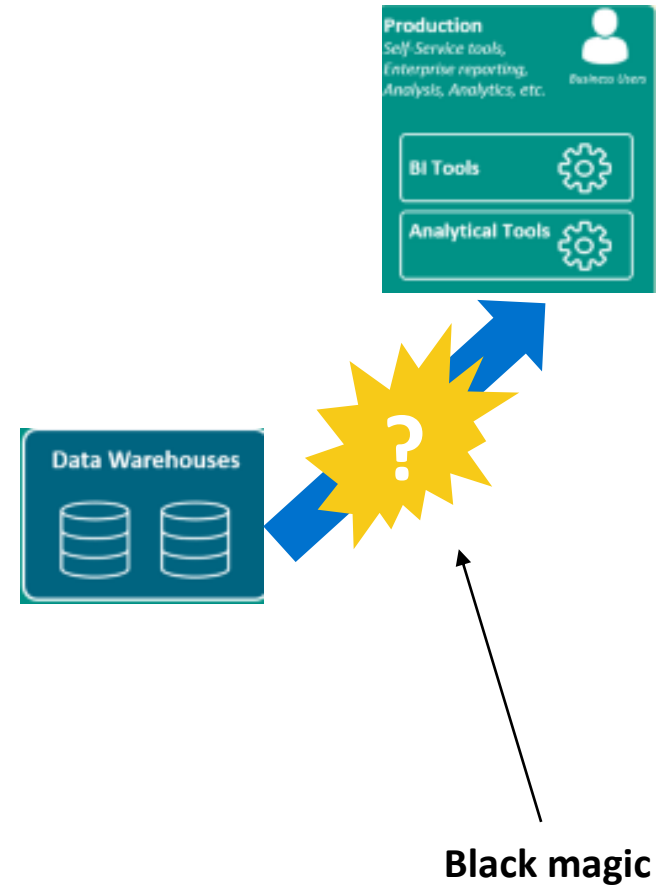# This is what "business" typically knows

**We have a cool branded data warehouse solution!**
**We can conquer the world!**
- Integrates several functionalities
- Promises to solve all the problems in one
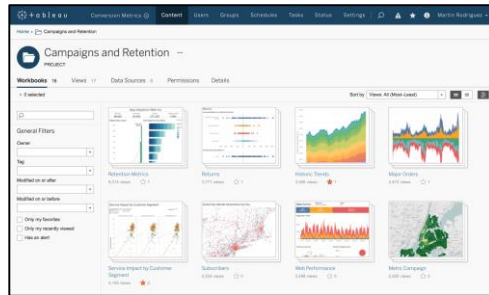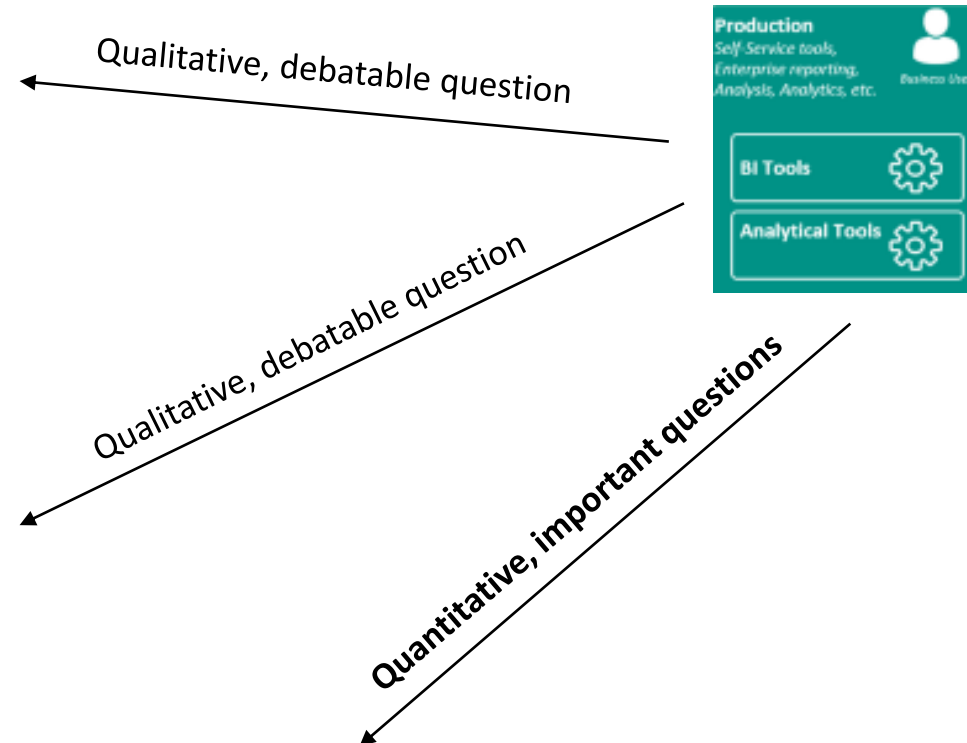- Costs a lot (this is why business knows!) ☺



Data Warehouses

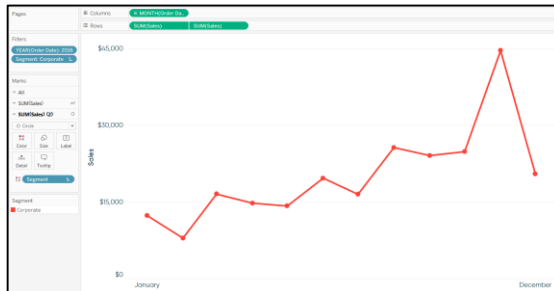# This is what "business" sees



**Black magic**

# This is what "business" asks

This graph is too complicated!



This graph is too simple!



Qualitative, debatable question

Qualitative, debatable question

**Quantitative, important questions**



- **Which version of the data did you use?**
- **Did you fix the upstream bug?**
- **Is the data updated?**
- **Is the data legally approved?**
- **Can we leverage more data?**
- **Can we re-use data from other projects?**
- **Who else is using this data?**

# Good luck in finding out… (still a simplified version!)
# How to make sense of such a complicated chain?



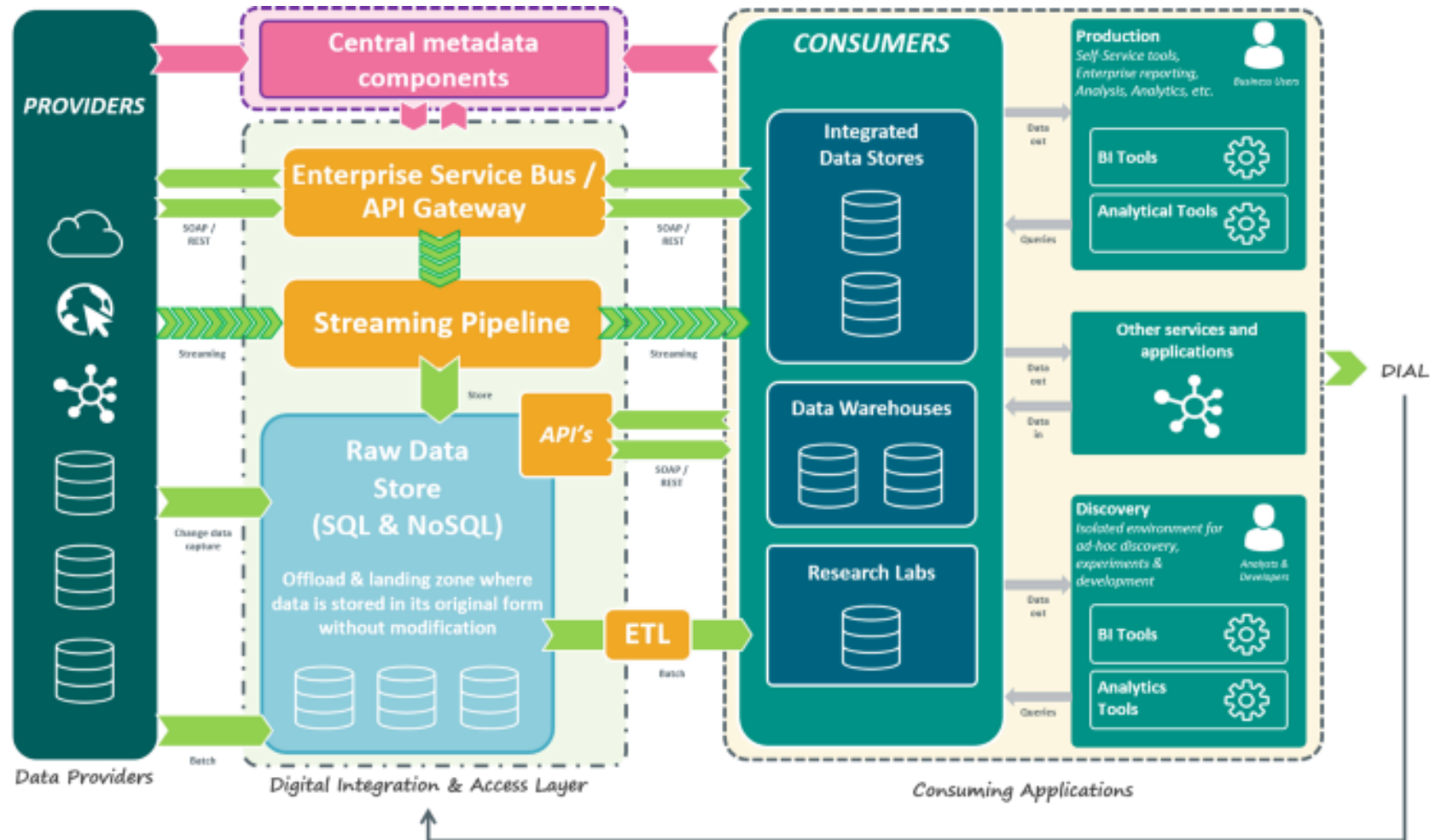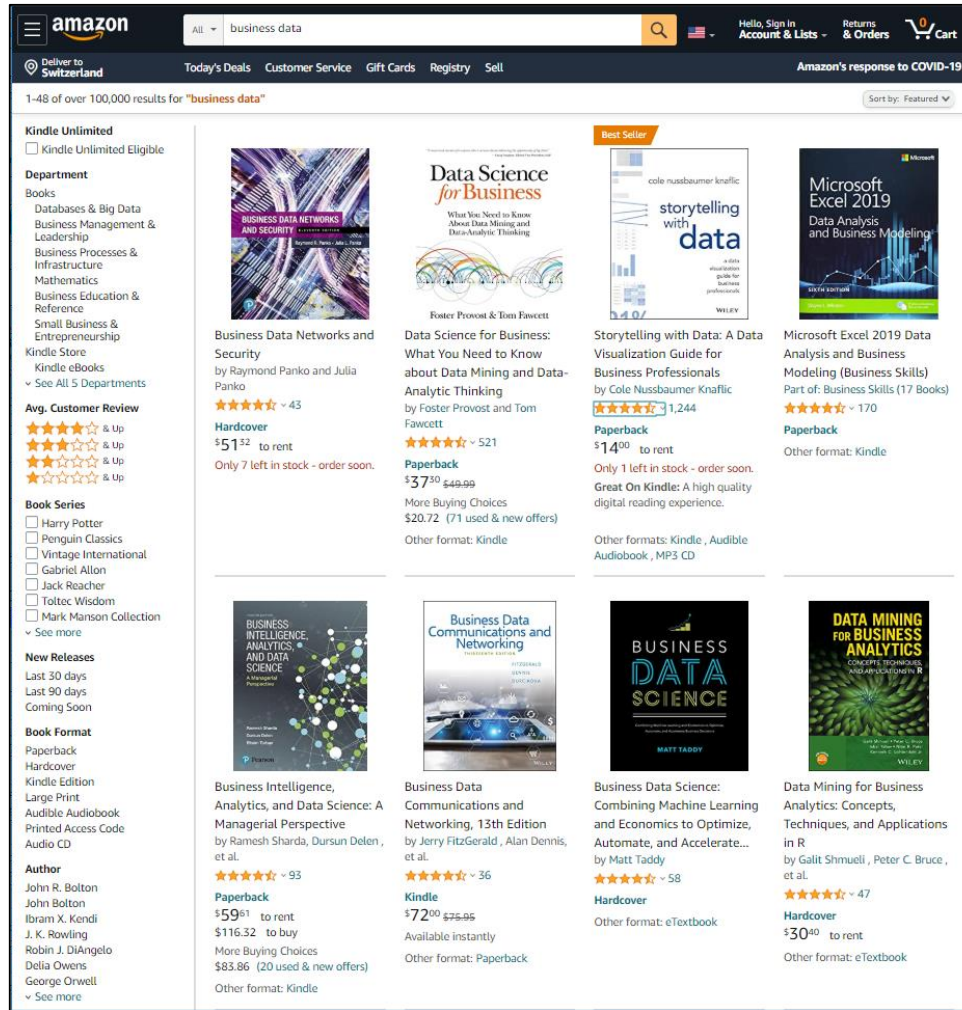Image source: https://medium.com/abn-amro-developer/abn-amros-data-integration-architecture-f33506a211c3

# Answer: use Metadata! (indicative, minimal list)

| Technical | Business | Operational (monitoring purpose) |
|---|---|---|
| • **Data source properties**<br>  • Platform (CH, R.O.W.)<br>  • Sourcing team<br>  • Version<br>• **Environments (PROD, PTA)**<br>  • Type of feed (interface)<br>  • Landing folder<br>  • Staging folder<br>• **Data Warehouse attributes**<br>  • Folder<br>• **Batch processor**<br>  • Cron schedule<br>  • Batch group<br>  • Batch timeout<br>• **Landing file format**<br>• **Staging file format**<br>• **Data source checks**<br>  • Check definition<br>• **Datasets properties**<br>  • Name<br>  • File pattern<br>  • Obfuscation details<br>  • Ingestion type<br>  • Versions<br>  • Dataset checks<br>    • Check definitions<br>  • Attributes<br>    • Name<br>    • Type | • **Data source**<br>  • Owner ID<br>  • Upstream application<br>  • Description<br>• **Initial requester**<br>  • Project<br>  • Business owner<br>• **Data category (GDPR)**<br>• **Service Level Agreement (SLA)**<br>  • Periodicity<br>  • Processing time<br>  • Reaction time<br>• **Data scope**<br>  • Division<br>  • Location<br>• **Originating/Sending/Receiving legal entities**<br>• **Delivery type**<br>• **Approvals**<br>  • Approval type<br>  • Approver<br>  • Date<br>• **Contacts**<br>  • Contact type<br>  • Contact name<br>  • Email<br>• **Data protection**<br>• **Datasets**<br>  • CID Class<br>  • Attributes | • **Batch name**<br>• **Execution id**<br>• **Execution start**<br>• **Execution end**<br>• **Processing Stages**<br>  • Stage start<br>  • Stage end<br>  • Stage result<br>  • Stage error message<br>• **Statistics**<br>  • Files statistics<br>  • Data statistics<br>• **File controls**<br>  • Control type<br>  • Control timestamp<br>  • Control result<br>  • Error message<br>  • Execution stage<br>• **Datasets**<br>  • Dataset name<br>  • Controls<br>    • Control type<br>    • Control timestamp<br>    • Control result<br>    • Error message |

# Dreaming about (meta)data representation…



- (meta)Data is perfectly **categorized**

- Technical, Business and Operational (meta)data **harmonized** and stored in a single place

- (meta)Data is **searchable**

- (meta)Data is available in **different formats**

- **Relations/similarities** between (meta)data are captured (more in the next slides)

- **Detailed descriptions** available to understand data

- **Quality** of the data is captured

How can we reach a similar level of representation?
**A good data management model and corresponding metadata collection is the key… and many (long forgotten) tools are available on the market!**

# All this might seem quite boring, so let's focus on a fun aspect: Data Lineage

- Data Lineage is a way of representing relations between datasets making use of graph analysis

- Graphs come with nodes (letters in the sketches below) and edges (lines in the sketches)
- Graphs can have different edge flavors (undirected, directed, weighted)



| Undirected | Directed | Weighted |
|---|---|---|
| Social media platforms (Facebook, Twitter, Linkedin, …) | Data platforms, Batch processing systems, Cooking (!) | GPS navigation systems (Google maps, Tom Tom, …) |

- Data platform graphs are typically «**Directed**» as we have **parent-child relations between nodes**

# Typical data warehouse lineage: end of the story?



- **Data warehouse tools typically provide lineage tools** (this example is taken from Palantir Foundry)

# Typical data warehouse lineage: what else do we need?



Typical main problems:
- Very restrictive access to metadata
- Not annotated **(i.e. no business logic!)**
- Limited graph analysis embedded functionalities

# Annotated lineage: technical solution

- Technical solution: use **networkx** (alternative solution: Java-based Neo4j)
  - An **open source Python package** for the creation, manipulation, and study of the structure, dynamics, and functions of complex networks.

# Annotated lineage: know your platform

- **Create and annotate** platform graph from data warehouse API and all the available metadata
  - **Solve ambiguities** and identify documentation errors
  - **Deal with different entities arising from business logic**: datasources, projects, others
  - **Cleanup** and add missing documentation
  - **Report** possible policy violations
  - **Answer business questions** in a fast and effective way

# Data warehouse lineage: typical facts and figures

- Typical Data warehouse lineage features:
  - O(100)   datasources
  - O(10)     projects
  - O(10k)   nodes
  - O(100k) direct connections
  - O(10M)  indirect connections

NB: each platform (R.O.W., CH)
and environment (UAT, PROD) has its
own distinct lineage graph!



Image source: https://blog.grakn.ai/life-on-the-edges-september-2016-cb4b900dd0a0
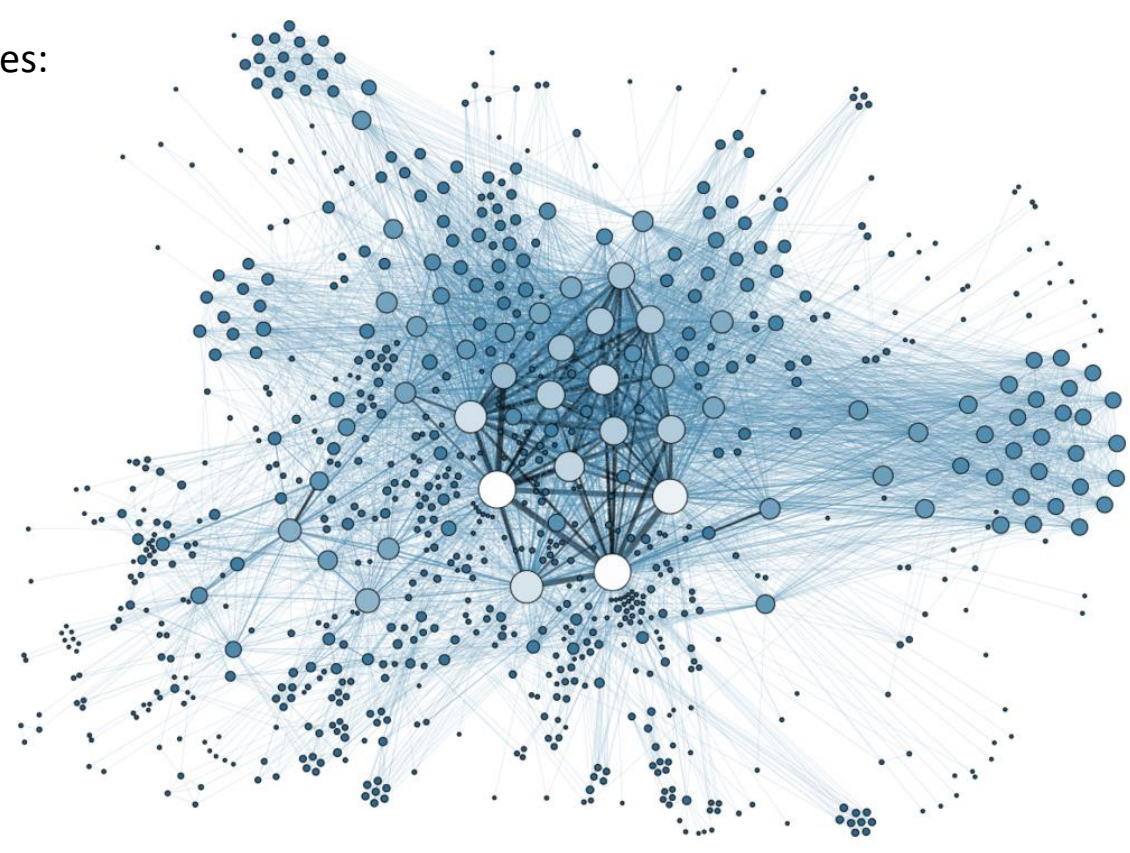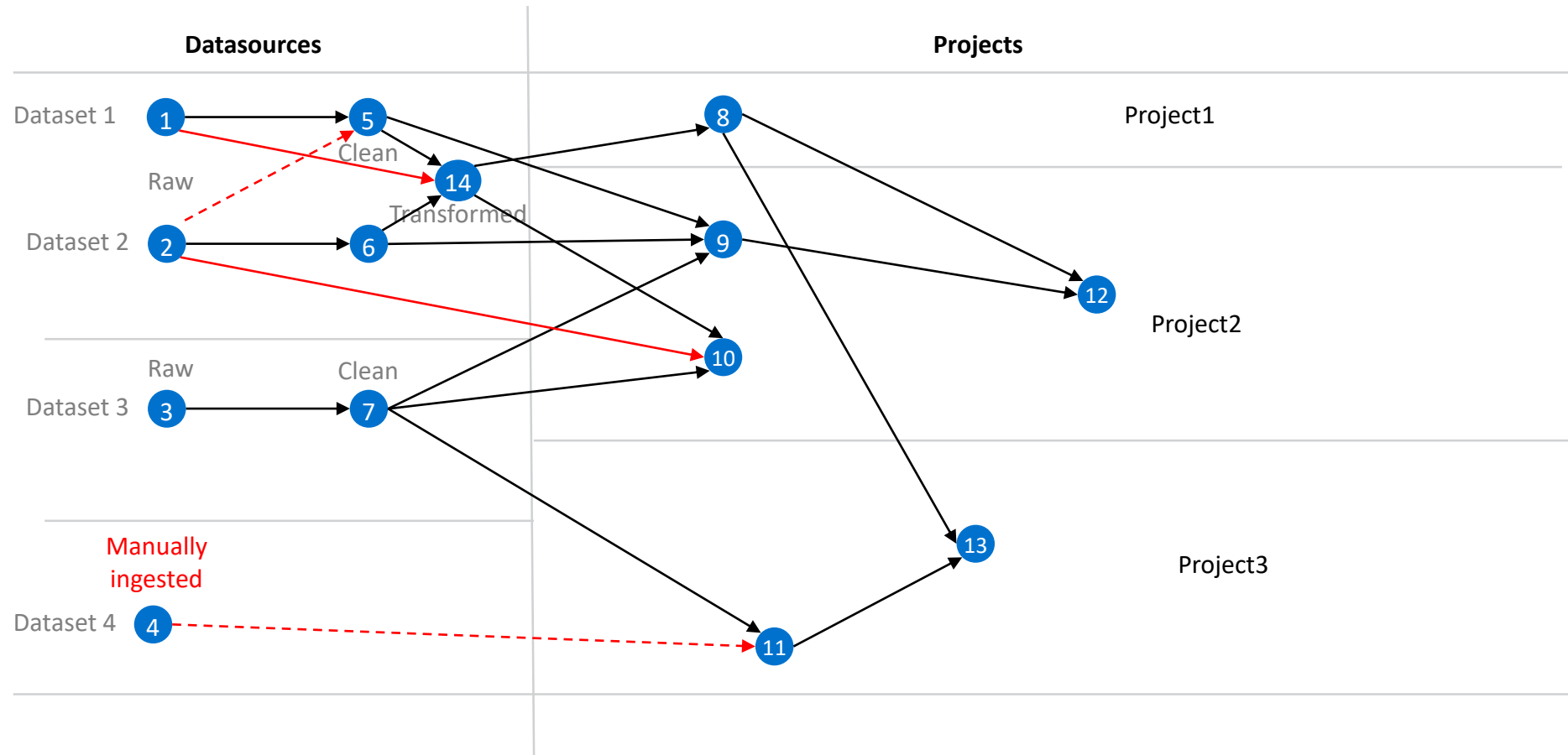
# Lineage simplified (!) schematics



- Two possible dependency relations for each node:
  - **Forward depencencies**: which other nodes depends on the current node?
  - **Backward depencencies**: which other nodes is the current node depending on?

# How to answer Lineage related business questions?

**Main ideas implemented to cope with system limitations:**

- Create a **«flattened» static representation of the graph using python package pandas** at the expense of duplicating information
  - One entry for each pair of connected nodes, either directly or indirectly
  - Updated daily

O(100k) rows

| Parent dataset | Child dataset | Direct connection |
|:---:|:---:|:---:|
| A | B | TRUE |
| B | C | TRUE |
| A | D | TRUE |

O(10M) rows

| Parent dataset | Child dataset | Direct connection |
|:---:|:---:|:---:|
| A | B | TRUE |
| **A** | **C** | **FALSE** |
| B | C | TRUE |
| A | D | TRUE |

- Indentify most common cases and **build reports** in excel or similar format for easy consumption from stakeholders
  - Provide overview and detailed versions

- Collect **lineage** for each platform (R.O.W., CH) and environment (UAT, PROD) **in one place**
  - Use **sharepoint** or similar technology to collect reports, so they can be consulted and searched online in the browser, and can be downloaded on the laptop for more refined analysis

# Lineage related business questions and reports

Typical business questions:

1. **Which projects depends on a given datasource?**
   - Report: list dependent projects for each datasource
2. **Business question: which other projects depend on a given project?**
   - Report: for each project, list other projects dependending on it
3. **Which datasources and other projects is a given project depending on?**
   - Report: list datasources and other projects dependencies for each project

Policy violations reports:

1. **Ingested but unused datasets**
2. **Stale datasources without ingestions in the last 3 months**
3. Raw datasets directly consumed by projects
4. Direct ingestion of datasources into project folders
5. Manually ingested datasets
6. Datasource depending on other datasources
7. Datasources using non-standard naming conventions

…

# Lineage and networkx DEMO

The material used in this demo can be found here:
https://github.com/perrozzi/networkx_example

# Conclusions

- Nothing new under the sun, but worth repeating ☺

- Data warehouses solve a whole lot of technical problems

- As usual, the **real problems** arise when the **technical solutions and operations need to co-exist with the business** expectations, evolving (conflicting) requirements and top-down decisions

- Data management and metadata are typically the most mistreated and disregarded information in this context
    - Still without them the technical solution will solve some problems and create (many) others

- Data management requires the coordination of several interested parties, both from the technical and business side

- Graph analysis and lineage allow an in-depth knowledge of data platform and can be extremely useful to answer business questions
    - And graph analysis is fun to implement and play with!