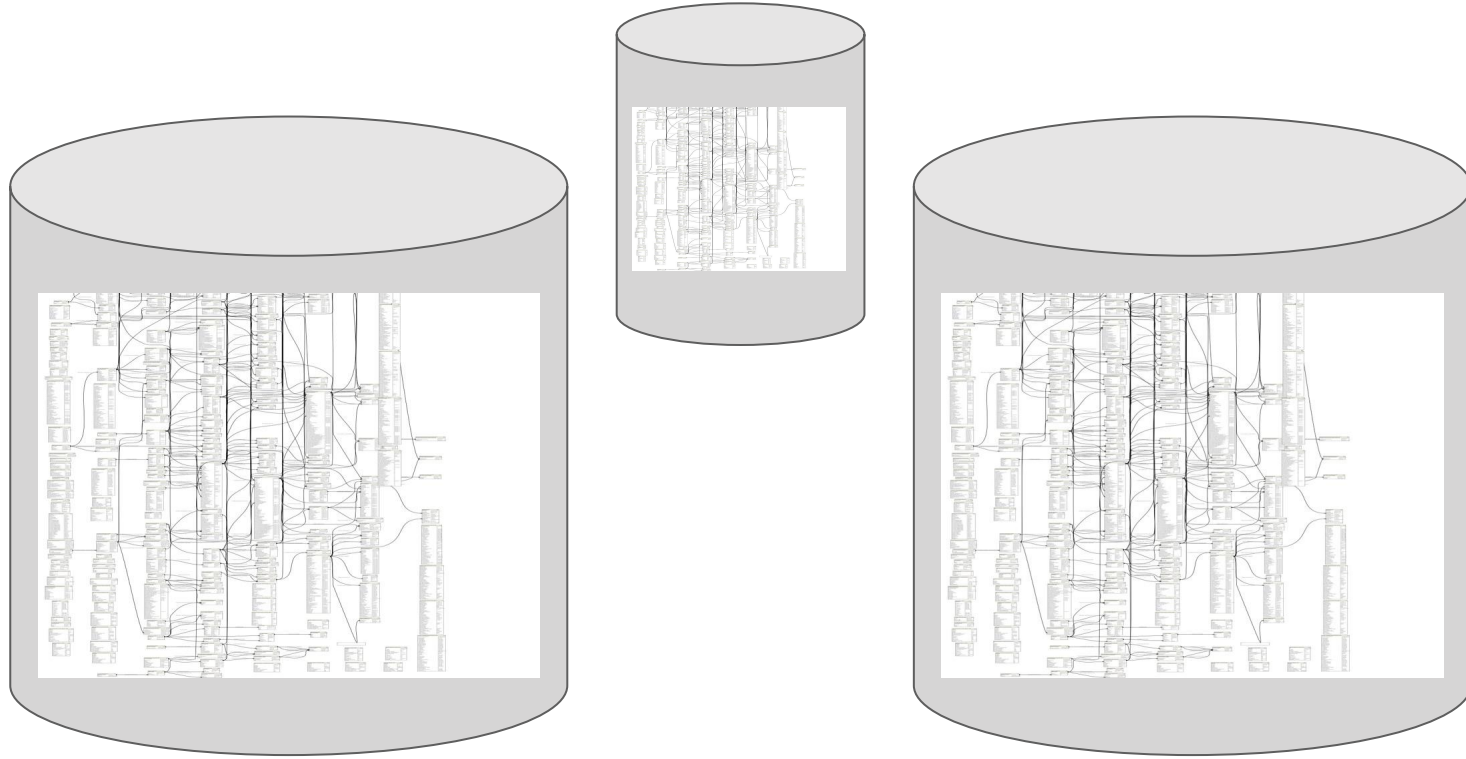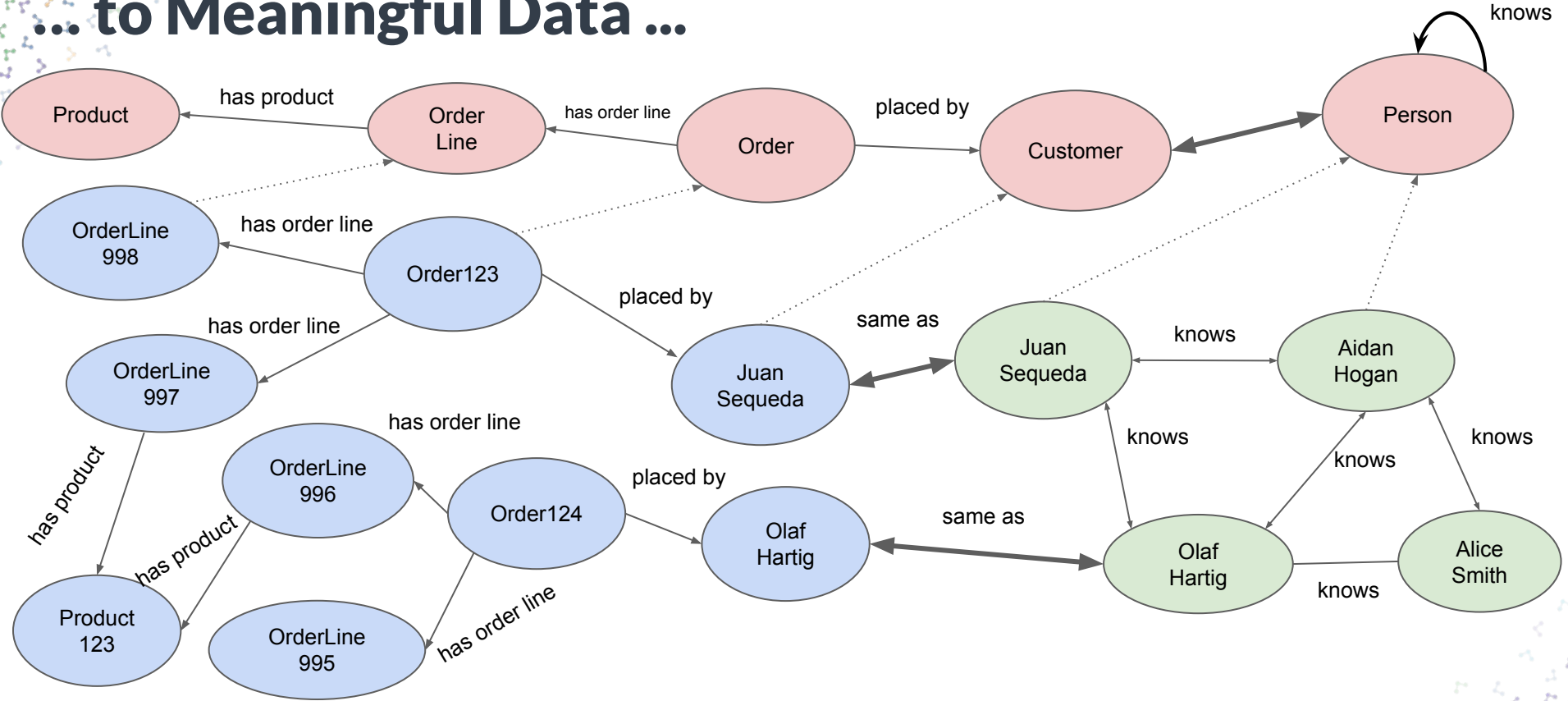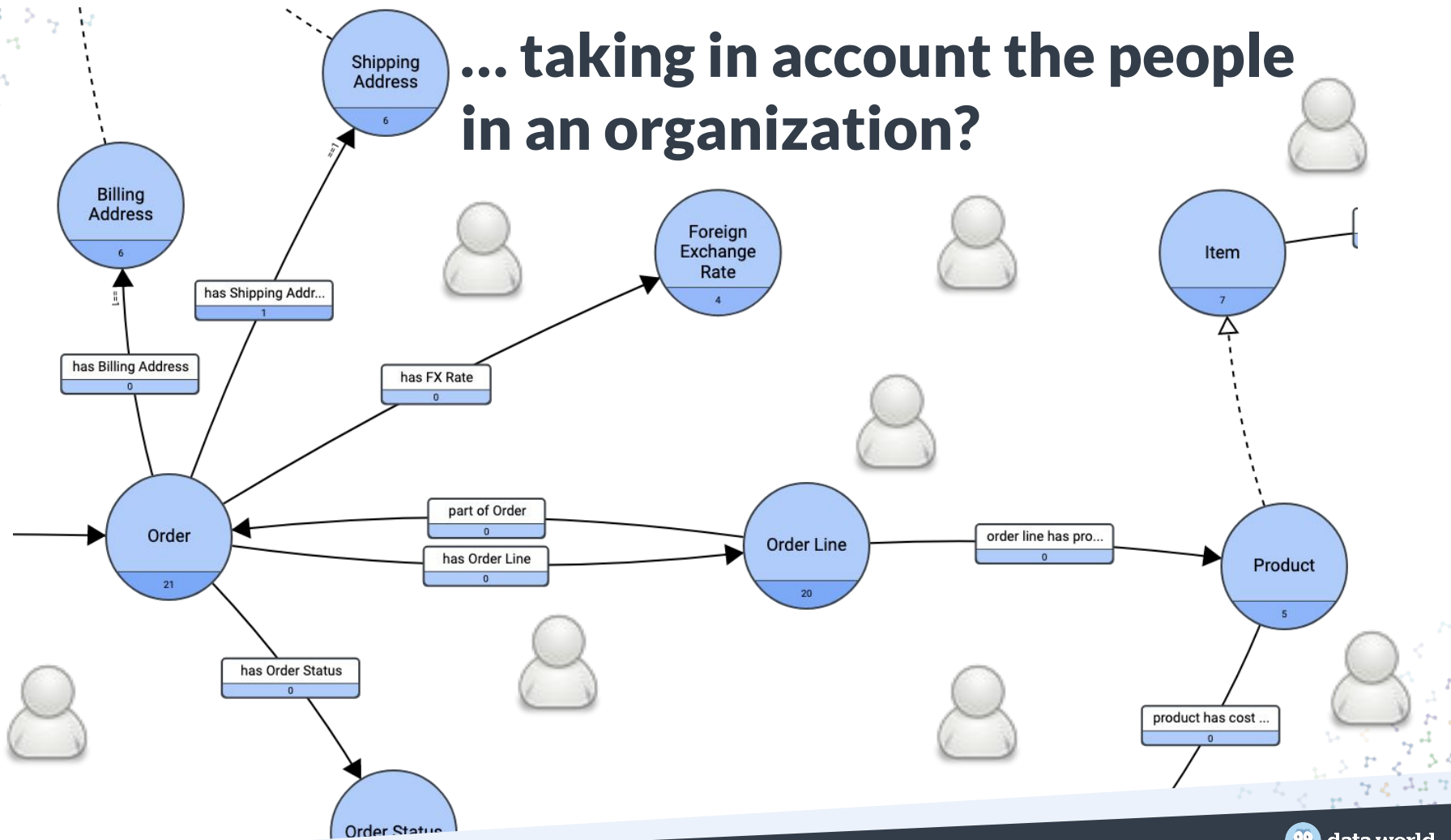# How do we get from Complex & Disparate Databases ...

# ... to Meaningful Data ...

... taking in account the people in an organization?

# The Socio-Technical Phenomena of Data Integration

data.world

**Juan Sequeda, PhD**
**Principal Scientist**

**@juansequeda**

data.world

# Takeaway



**My Thesis:**
We have been studying the phenomena of data integration from a technical perspective ...

**My Question:**
*How can we best combine people and technology to improve data integration?*

IT

Data Analyst

Data Scientist

Data Steward

ETL Developer

Subject Matter Expert

Data Engineer

Business User

... while ignoring the **social** aspect

data.world

# Data Integration

WHY?

*"Data integration is the problem of combining data residing at different sources, and providing the user with a unified view of these data."*

- Schema/Ontology Matching
- Record Linkage/Instance Matching
- Incomplete Data
- Data Quality
- Materialized vs Virtual
- ...

Query

**Global schema**

**Mapping**

| $R_1$ | $C_1$ | $D_1$ | $T_1$ |
|-------|-------|-------|-------|
|       | $c_1$ | $d_1$ | $t_1$ |
|       | $c_2$ | $d_2$ | $t_2$ |

**Source schema**

**Source schema**

Source 1

Source 2

Maurizio Lenzerini (2002). "Data Integration: A Theoretical Perspective" . PODS 2002.

data.world

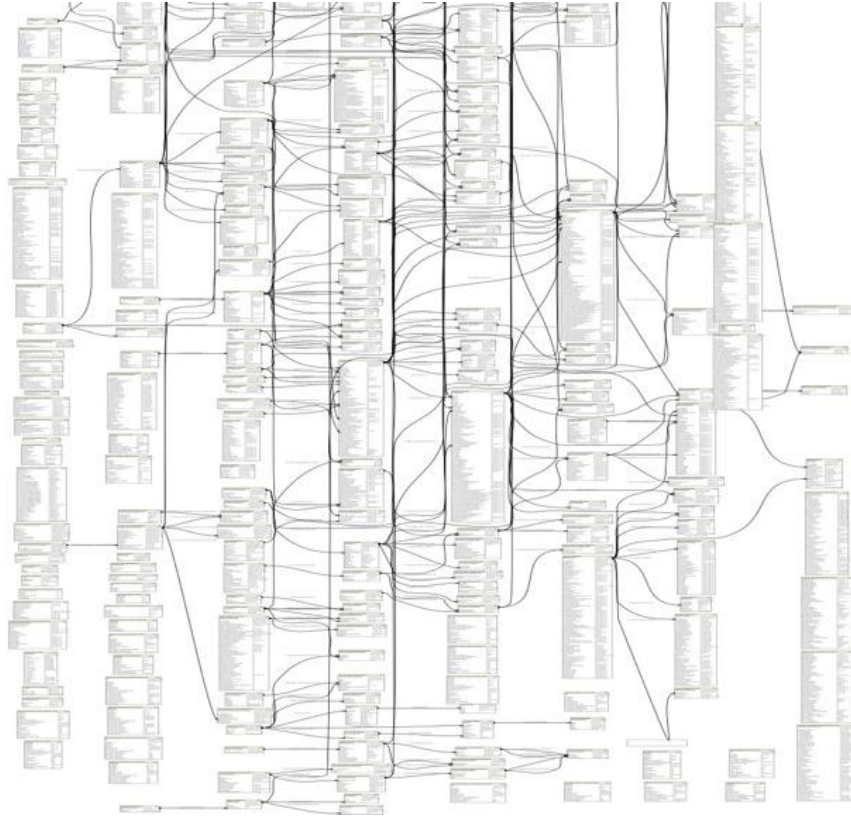# Real World Socio-**Technical** Problems

data.world

# Technical 1: Understand the Source

**Too many tables and attributes**

**Complex Relationships**

**Data experts unavailable**

**Master databases are off limits**

**Impossible to understand naming**

**Data is application centric**

**Documentation non-existent**

**Data quality unknown**

data.world

# Technical 2: Understand the Target

Sophisticated use of competency questions, test-driven development, ontology design patterns, reuse.

Populating ontology with data coming from a relational database is an afterthought and not part of methodologies.

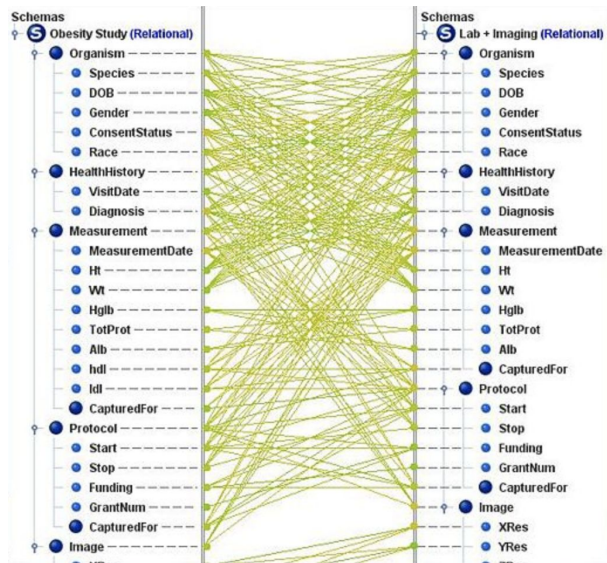Numerous upper level ontologies for reuse: Good Relations, FIBO, Gist, Schema.org, etc.

data.world

# Technical 3: Understand the Mappings

**Ontology/Schema matching between relational schemas and ontologies works … in theory. But not in practice.**

MasterOrder,
Order,
P_Order,
Order_P

segment1,
segment2,
…,
segment99



1-1
correspondence
between
table-classes and
columns-properties
are <u>rare</u>

**Not plausible we will ever have large amounts of schemas & mappings to train ML models.**
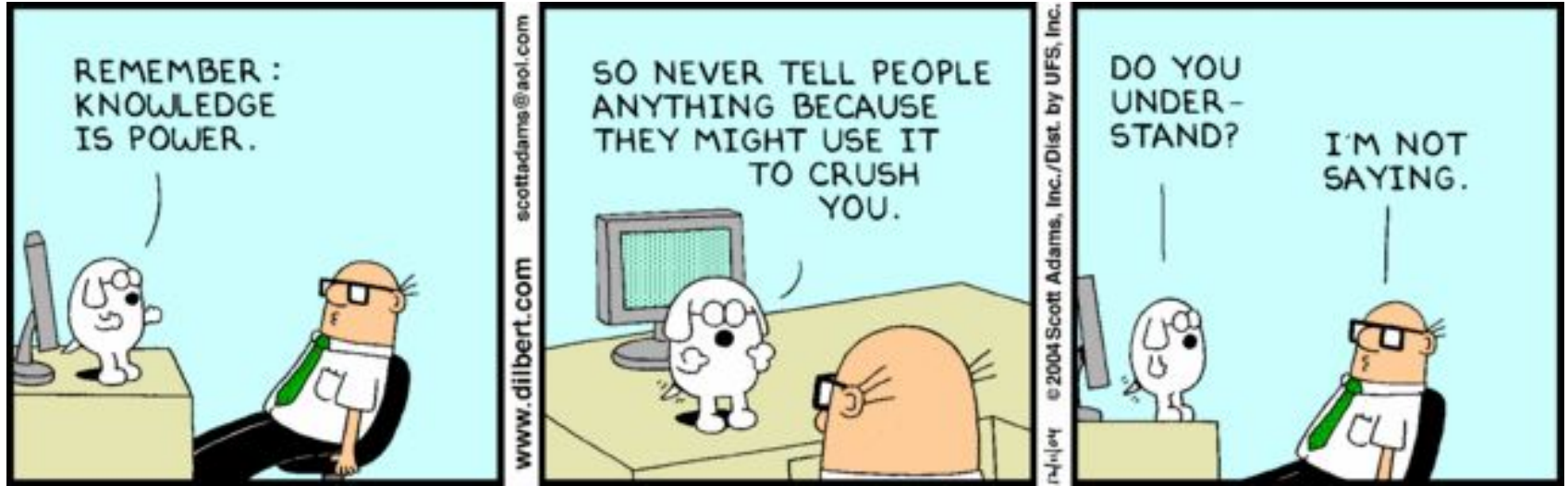
data.world

# Social 1: Knowledge Hoarding
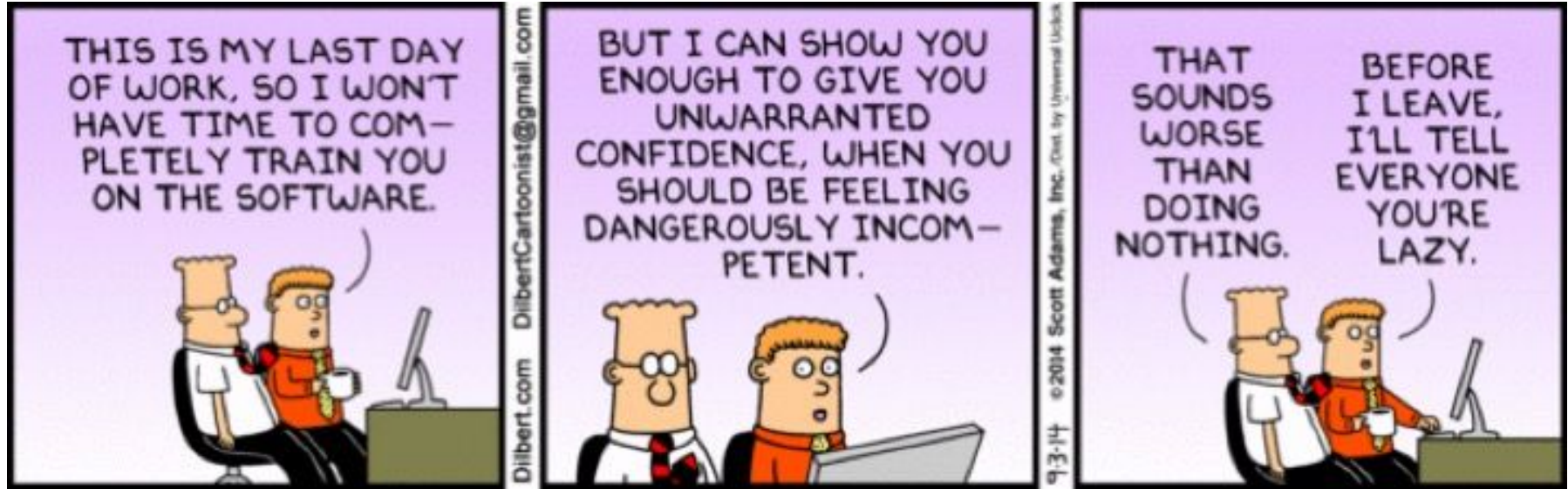
# Social 2: Knowledge Aspiration

# Social 3: Knowledge Retirement

# Real World **Socio-Technical** Problems

E-commerce: 118,595

Shipping: 114,324

Finance: 116,211

# Business Question

**How many orders were placed in December 2019?**

data.world

**E-commerce: 118,595**

*When the user clicks "Order" on the website.*

**Shipping: 114,324**

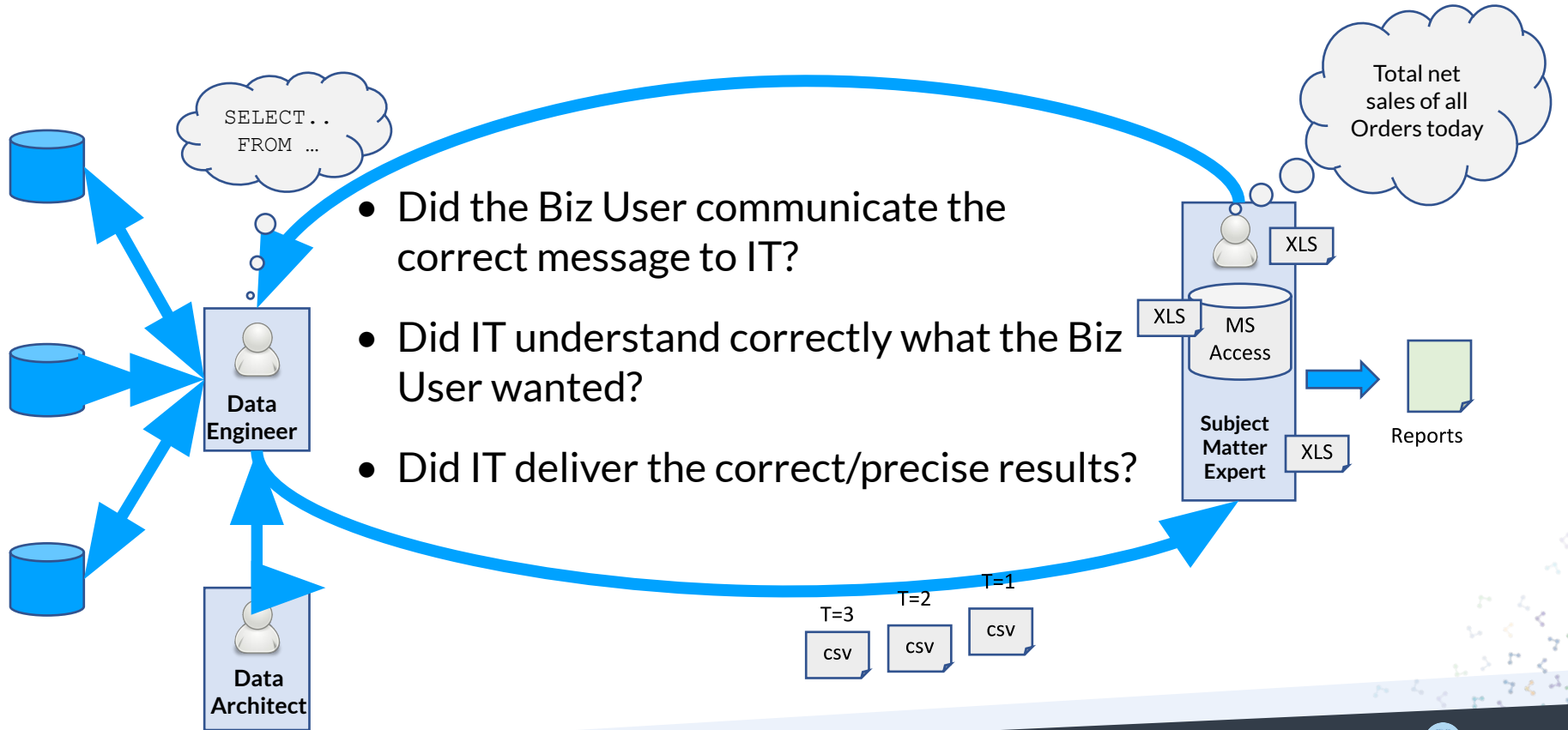*When the customer has received the product*

**Finance: 116,211**

*When it comes out of the billing system and the CC has been charged*
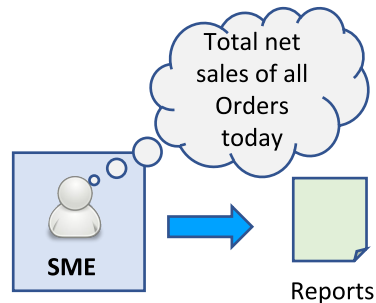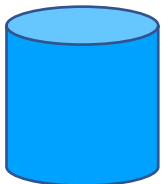
**What do you mean by ...**

What is an **Order**?

data.world

# Spreadsheet Approach



SELECT.. FROM …

Total net sales of all Orders today

- Did the Biz User communicate the correct message to IT?

- Did IT understand correctly what the Biz User wanted?

- Did IT deliver the correct/precise results?

Data Engineer

Data Architect

XLS

XLS

MS Access

XLS

Subject Matter Expert

Reports

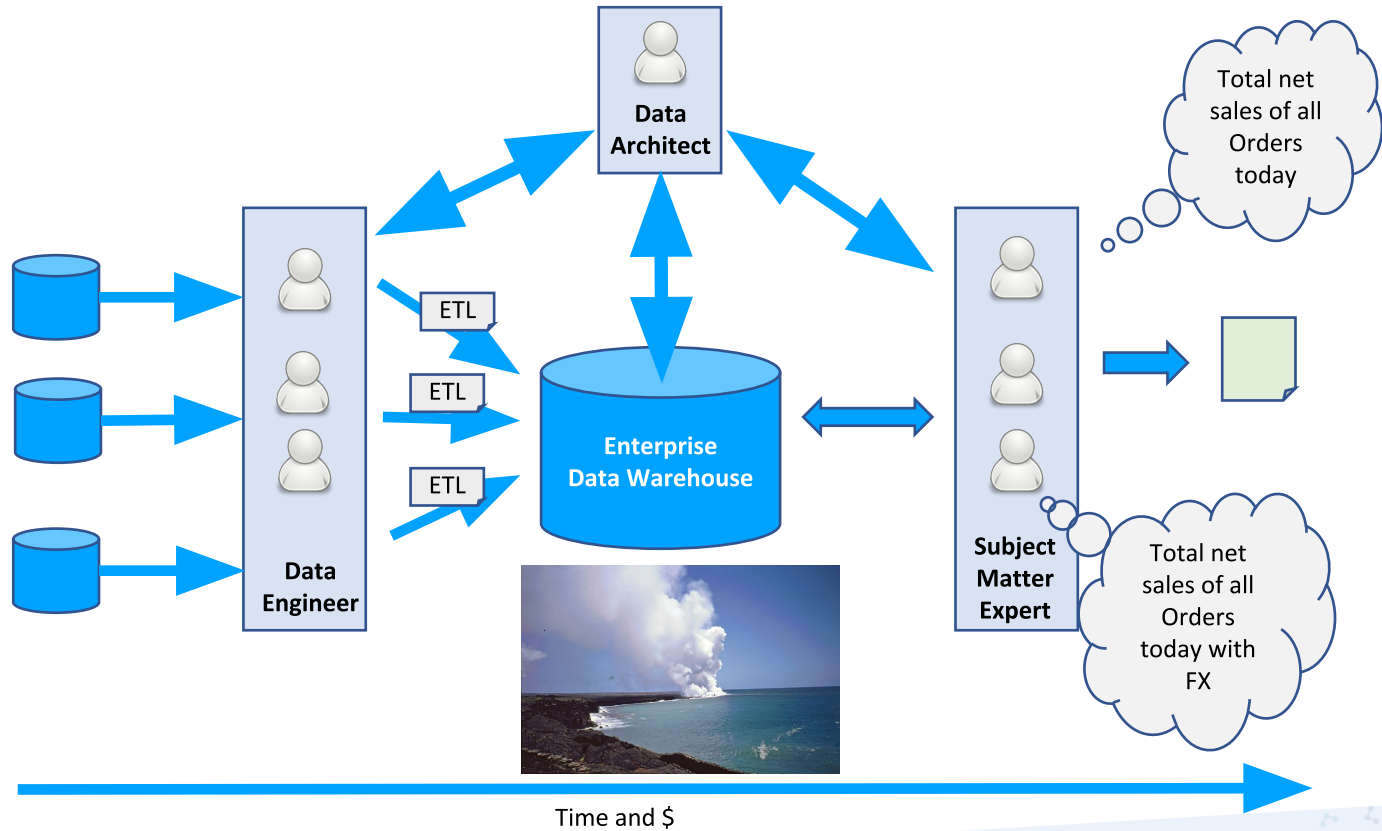T=3 csv

T=2 csv

T=1 csv

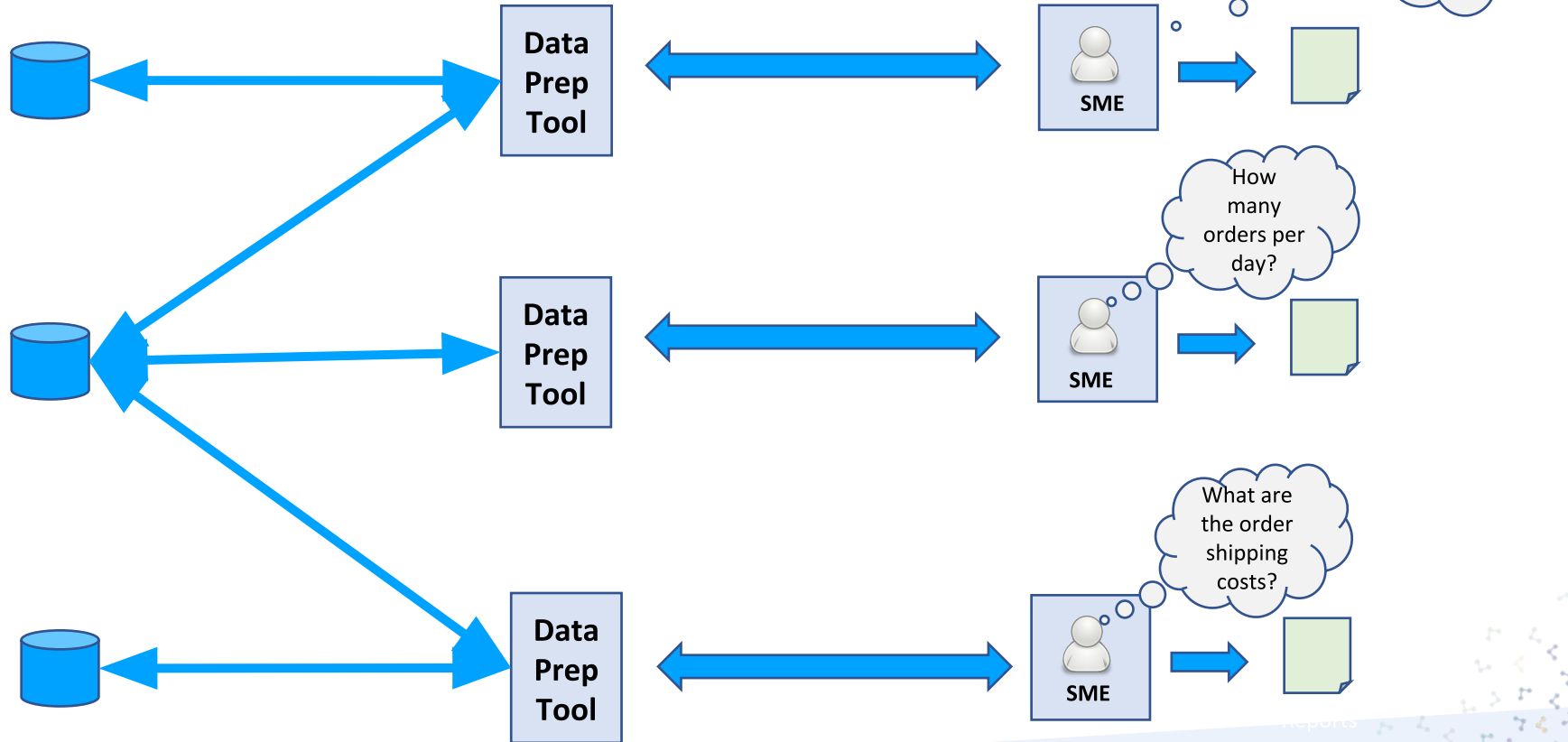data.world

# Query Approach



```
isnull(POE.dbo.OrderSku.PRICE,0)*isnull(POE.dbo.OrderSku.QUANTITY,0)-isnull(POE.dbo.OrderSku.TOTALDISCOUNTAMOUNT,0)-isnull(POE.dbo.Or

CASE WHEN POE.dbo.POEMasterOrders.MOOrderTypeID like 'REP%' THEN 0 ELSE
                isnull(POE.dbo.OrderSku.OverrideCVEarned,0)END As Total_RV_OverrideCVEarned,

    CASE WHEN POE.dbo.[Order].Subtotal >0 THEN CASE WHEN ACM.dbo.AppBeeClass.ClassID = 'PreferredCustomer' THEN 0.15*isnull(POE.dbo.Or
        WHEN ACM.dbo.AppBeeClass.ClassID = 'Customer' THEN 0.35*isnull(POE.dbo.OrderSku.CommissionableValueEarned,0)
        ELSE 0 END END AS Comm,

            CASE WHEN isnull(POE.dbo.OrderSku.PRICE,0)*isnull(POE.dbo.OrderSku.QUANTITY,0)-isnull(POE.dbo.OrderSku.TOTALDISCOUNTAMOU
        = 0 THEN 0 ELSE
                isnull(POE.dbo.OrderSku.PRICE,0)*isnull(POE.dbo.OrderSku.QUANTITY,0)/1+
        (isnull(POE.dbo.OrderSku.TaxAmountIncluded,0)/
        (isnull(POE.dbo.OrderSku.PRICE,0)*isnull(POE.dbo.OrderSku.QUANTITY,0)-isnull(POE.dbo.OrderSku.TOTALDISCOUNTAMOUNT,0)-isnull(POE.
                AS SRP_Est_BeforeTax,

Sub_fxr.To_USD

FROM
 POE.dbo.OrderSku LEFT OUTER JOIN POE.dbo.PoeProductTypeCode ON (POE.dbo.OrderSku.PRODUCTTYPECODEID =POE.dbo.PoeProductTypeCode.PoePr
 LEFT OUTER JOIN POE.dbo.[Order] ON (POE.dbo.OrderSku.ORDERID=POE.dbo.[Order].OrderId)
 LEFT OUTER JOIN POE.dbo.POEOrders ON (POE.dbo.OrderSku.ORDERID=POE.dbo.POEOrders.OrderID)
 LEFT OUTER JOIN POE.dbo.POEMasterOrders ON (POE.dbo.POEOrders.MasterOrderID=POE.dbo.POEMasterOrders.MasterOrderID)
 LEFT OUTER JOIN POE.dbo.Sku on (POE.dbo.OrderSku.SKUID= POE.dbo.Sku.SkuID)
 LEFT OUTER JOIN POE.dbo.OrderStatus ON (POE.dbo.[Order].OrderStatusId=POE.dbo.OrderStatus.OrderStatusId)
 LEFT OUTER JOIN ACM.dbo.GENORDERTYPE ON (POE.dbo.POEMasterOrders.MOOrderTypeID = ACM.dbo.GENORDERTYPE.OrderTypeID)
 LEFT OUTER JOIN ACM.dbo.AppBeeClass on (POE.dbo.POEOrders.AppBeeClassGuid=ACM.dbo.AppBeeClass.AppBeeClassGuid)

 LEFT OUTER JOIN

(Select distinct
        POE.dbo.poeorderitemdisc.orderskuid,
POE.dbo.PoeDiscOption.Name As PlanDesc,
POE.dbo.poediscplan.Description As PlanName1

 from
POE.dbo.poeorderitemdisc
 join POE.dbo.PoeDiscPlanLevelDisc on POE.dbo.PoeDiscPlanLevelDisc.PoeDiscPlanLevelDiscId = POE.dbo.poeorderitemdisc.PoeDiscPlanLevelD
 join POE.dbo.PoeDiscOption on (POE.dbo.PoeDiscOption.poediscoptionid = POE.dbo.PoeDiscPlanLevelDisc.poediscoptionid and POE.dbo.PoeDi
 join POE.dbo.poediscplanlevel on poe.dbo.poediscplanleveldisc.poediscplanlevelid = poe.dbo.poediscplanlevel.poediscplanlevelid
 join POE.dbo.poediscplan on poe.dbo.poediscplanlevel.poediscplanid = poe.dbo.poediscplan.poediscplanid
 where
                            POE.dbo.PoeDiscOption.SystemKeyword in ('BuildToOrderDiscount','DiscountedSkuFromList'
```

Total net sales of all Orders today

SME

Reports

data.world

# Data Warehouse Approach

# Data Wrangling Approach

# Why is this hard?

data.world

# Strings to Things

**What is an Order?**

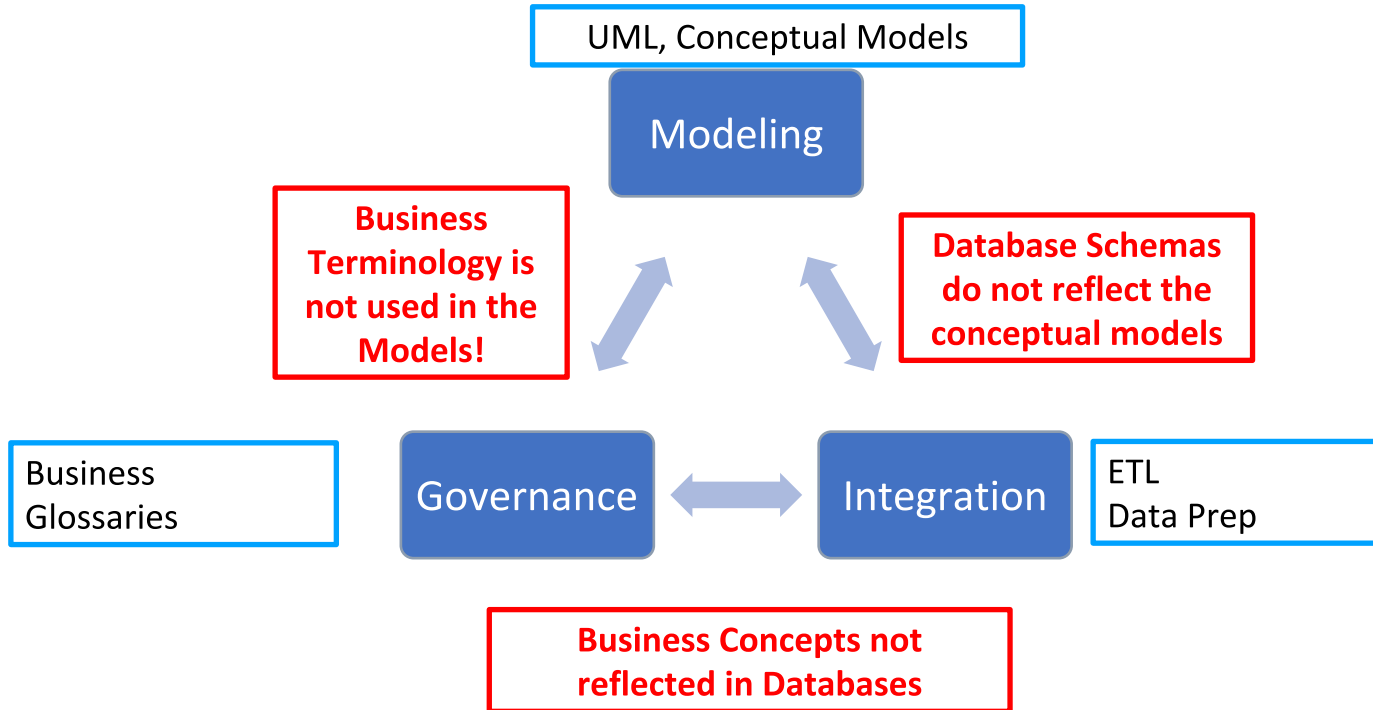An order is if it had shipped or the accounts receivable had been received.

```
SELECT moid
FROM masterorder m
JOIN order o on m.oid = o.oid
WHERE ordertype IN (2,3)
```

**What is the net sales of an order?**

Net sales of an order is calculated by subtracting the tax and the shipping cost from the final price and also adjusting based upon the discount given. However, if the currency of the order is not in USD or CAD, then the shipping tax must be subtracted.

```
SELECT
    moid,
    o.ordertotal - ot.finaltax -
    CASE WHEN o.currencyid in ("USD", "CAD")
        THEN o.shippingcost
        ELSE o.shippingcost = ot.shippingtax
    END as ordernetsales
FROM masterorder m
JOIN order o on m.oid = o.id
JOIN ordertax ot on o.oid = ot.oids
```

data.world

# People and Tool Fragmentation

# Real World Example

# An E-Commerce Case Study

## Background

- Customer invested in expensive EDW

- Business Users were skeptical of the data feeding the BI reports

- EDW was not being used.

- BU generating reports directly from sources.

## Challenge

- IT quickly became the bottleneck

- Friction between BU and IT due to lack of agreed terminology.

- Tribal knowledge in Excel, MS Access.

- Different answers for the same question.

- How is Tableau going to be successful?

data.world

# Customer Need

- **Consistent, understandable and trusted data view across the multiple relational databases.**

- **Tableau needs to consume the trusted data.**

- **Large number of business users should generate reports using the same trusted data.**

- **Agile approach in order to start showing value quickly.**

data.world

# THE TECHNOLOGY FALLACY

## HOW PEOPLE ARE THE REAL KEY TO DIGITAL TRANSFORMATION

GERALD C. KANE, ANH NGUYEN PHILLIPS,
JONATHAN R. COPULSKY, AND GARTH R. ANDRUS

*"The mistaken belief that just because business challenges are caused by digital technology, that they also need to be solved by digital technology."*

# People

## Data Engineer

Understand database schemas, including how the data are interconnected.

## Knowledge Scientist

Serves as the communication and developer bridge between Data Engineers and Business Users

Data Access

"Geeky Person"



Business Modeling

"People Person"

## Business User

SME who understand the business

data.world

# Knowledge Scientist vs Data Scientist

"Most data scientists spend only 20 percent of their time on actual data analysis and 80 percent of their time finding, cleaning, and reorganizing huge amounts of data, which is an inefficient data strategy"

Knowledge Scientist

Knowledge Scientist

Data Scientist

data.world

# What data do we have?

1) Catalog

Let's catalog the data!

data.world

3) Knowledge Graph
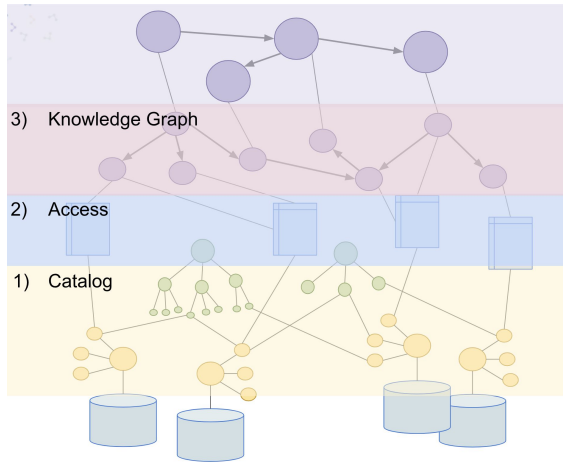
2) Access

1) Catalog

I don't understand my data!

Let's add knowledge to the data!

data.world

# Data Cloud

A data cloud is where your data is available to your people and your machines – it's where your data assets are leveraged.
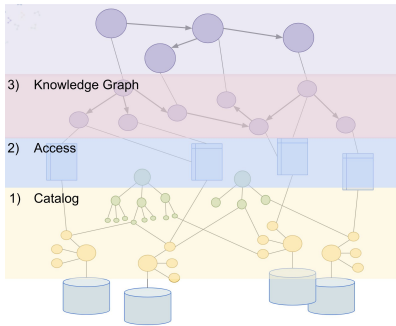


**+**



**=**

## Data Cloud

The 2019 LinkedIn Top Startups Are Growing Fast

1)     Snowflake
2)     …

3)   Knowledge Graph

2)   Access

1)   Catalog
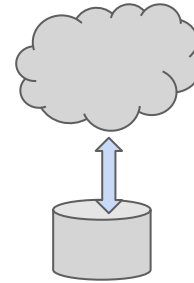
data.world

# Hybrid Data Cloud

The data cloud will evolve to be hybrid, combining elements of your on-prem systems and systems native to public clouds.
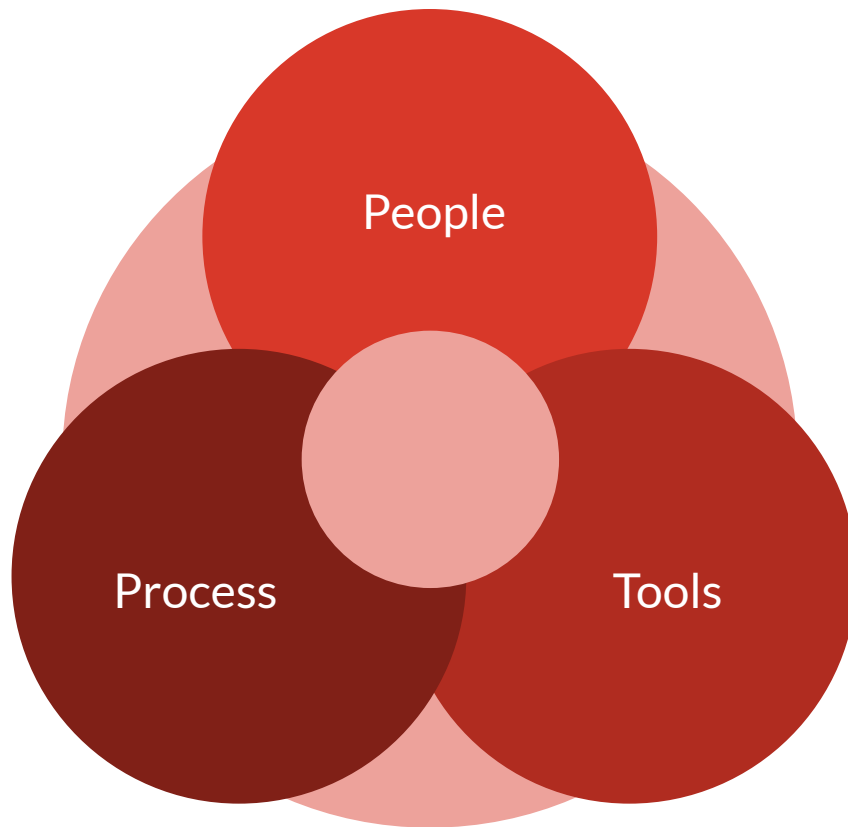


**+**



**+**



Knowledge Graph Virtualization

**=**

**Hybrid Data Cloud**
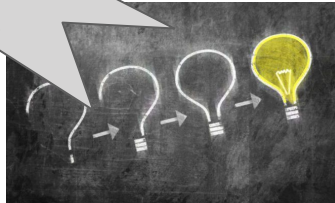
**Knowledge Report**

**Pay-as-you-go Methodology**

**Knowledge Capture**
1. Analyze as-is process
2. Collect Documentation
3. Develop Knowledge Report

**Business Question**

THIS IS HOW WE MEASURE SUCCESS!
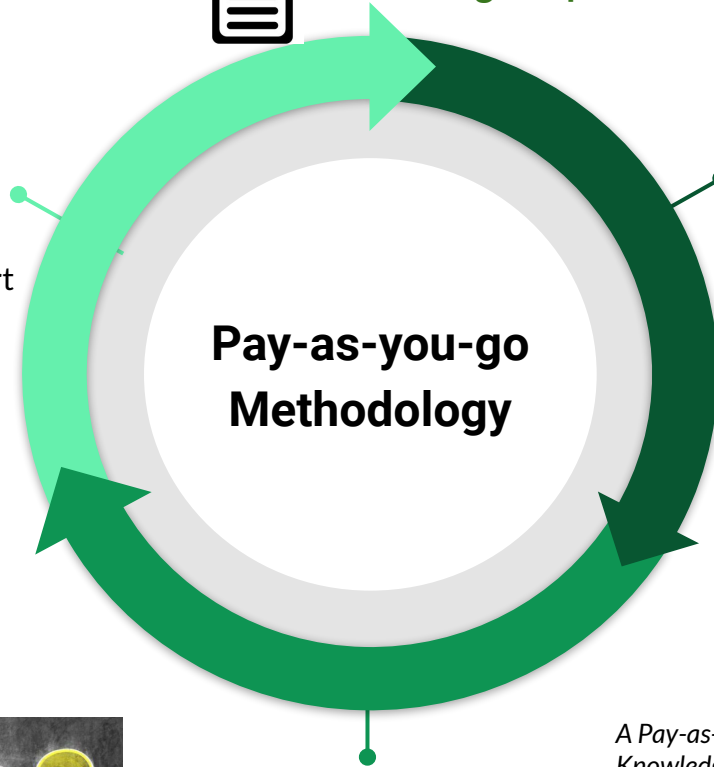
**Business Answer**

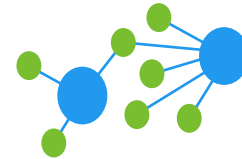**Knowledge Implementation**
4. Create/Extend Ontology
5. Implement Mapping
6. Create Extract Queries
7. Validate Data

**Enterprise Knowledge Graph**

*A Pay-as-you-go Methodology to Design and Build Enterprise Knowledge Graphs from Relational Databases. ISWC 2019*

**Self Service Analytics**
8. Build Report
9. Answer Business Question
10. Move to Production

data.world

# 1) Analyze as-is Process and Workflow

| | |
|---|---|
| **WHAT** is the business problem/question? | How many orders were placed in a given time period per their status? |
| **WHY** do we need to answer these questions? | Depending on whom is asked, different answers can be provided. Unaware of the source of the problem, the executives are vexed by inconsistencies across established business reports. |
| **WHO** produces/consumes the data? | The Finance department, specifically the CFO |
| **HOW** is this the business question answered today? | A business analyst asks the IT developer for this information every morning. |
| **WHERE** is the data? | There is a proprietary Order Management System and Oracle E-Business Suites. |
| **WHEN** will it be consumed? | Every morning they want to know this number |

# 2) Collect Documentation

- **Focus on HOW and WHAT**

- **Documentation, Wiki, SQL queries, Excel, ETL scripts, MS Access**

- **Interviews to understand the people and tech workflow**
  - **Who talks to who?**
  - **What is being shipped around?**
  - **Reverse engineer reports, queries, datasets**

# 3) Knowledge Report

**Data Engineer**

**Knowledge Scientist**

**Business User**

Identify which database schema elements contains related data

Understand the business questions, recognize key concepts and relationships, identify the terminology

Disagreements?
Focus on the business questions to drive consensus

Knowledge Report can be understood by everyone!

Knowledge Reports mimics the Intermediate Representations (IRs) from METHONTOLOGY

data.world

# 3) Knowledge Report: Concept



| Concept Name | Order |
|---|---|
| Concept ID | Order |
| Concept Instance ID | `moid` |
| Table Name/Query | `SELECT moid FROM masterorder m JOIN order o on m.oid = o.oid WHERE ordertype in (2,3)` |

# 3) Knowledge Report: Attribute



| | |
|---|---|
| **Attribute Name** | Order Date |
| **Attribute ID** | orderDate |
| **Applied to Concept** | Order |
| **Table Name/Query** | `select moid, orderdate from masterorder m join order o where m.oid = o.id` |
| **Column Name** | `orderdate` |
| **Datatype** | `date` |
| **Is NULL possible?** | No |
| **Cardinality** | 1:13) Knowledge Report: Concept |

# 3) Knowledge Report: Relationship



| | |
|---|---|
| **Relationship Name** | has order status |
| **Relationship ID** | hasOrderStatus |
| **From Concept** | Order |
| **To Concept** | Order Status |
| **Table Name/Query** | `select moid, ostid, max(orderstatusdate) from OrderStatus group by orderstatusdate` |
| **Cardinality** | 1:1 |

# 3) Knowledge Report: Tabular Extract



```
SELECT ?Order_Number ?Order_Date ?Order_Status
WHERE {
 ?o a :Order;
     :orderNumber ?Order_Number;
     :orderDate ?Order_Date;
     :hasOrderStatus [
          :orderStatus ?Order_Status
     ].
}
```

| Order Number | Order Date | Order Status |
|---|---|---|
|  |  |  |
|  |  |  |

**Knowledge Capture**  **Knowledge Implementation**

# 4) Create/Extend Ontology



**Knowledge Report**

```
ec:Order rdf:type owl:Class ; rdfs:label
"Order" .
ec:orderDate rdf:type owl:DatatypeP
    domain ec
     nge xsc            rdfs:l
    te" .
     rderStatus rdf:type
owl:objectProperty ; rdfs:label "has order
status"
   rdfs:domain ec:Order ; rdfs:range
ec:OrderStatus .
```

**OWL Ontology** *

**Create the ontology using Gra.fo**

\* Property Graph Schema too

data.world

# 5) Implement Mapping

**Knowledge Report**

```
map:m1 a rr:TriplesMap ;
  rr:logicalTable  [ rr:sqlQuery  "select
moid from masterorder m join
                          order o
             oid whe      e in  2
         ectMap      [ rr:class
          emplate
"http://www.e-commerce.com/data/Order/{moid
} ] .
```

Order

has order status
0

Order Status
0

**R2RML Mapping**

**Implement the mappings using Gra.fo**

# 6) Extract/Tabular Queries

**Knowledge Report**

```
SELECT ?Order_Number ?Order_Date ?Order_Status
WHERE {
?x a :Order;
  :orderNumber ?Order_Number;
  :orderDate ?Order_Date;
  :hasOrderStatus [
    :orderStatusName ?Order_Status;
  ]
}
```

**SPARQL Query**

**Execute the SPARQL query on Ultrawrap ETL/NoETL, which uses the R2RML mapping from the previous step.**

data.world

# 7) Validate Data

- **Counts: compare the number of results between extract and source**

- **Null: checking the validity of NULL values**

- **Duplicates: Check that expected Cardinality holds**

- **Sharing sample data to business users in a spreadsheet**

- **Creating sample visualizations in a BI tool**

data.world

**Knowledge Implementation**

**Self Service Analytics**

data.world

# 8) Build Business Report

- **The ontology enables the simplified view but at the end, it is the tabular extract that the business user wants to access.**

| Order Number | Order Date | Order Status |
|---|---|---|
|  |  |  |
|  |  |  |

😊

```
SELECT * FROM Orders
```

☹️

```
SELECT m.moid as OrderNumber,
       o.orderdate as OrderDate,
       ost.statustype as OrderStatusName
FROM masterorder m JOIN order o ON m.oid = o.oid
JOIN (SELECT moid, ostid, max(orderstatusdate)
      FROM OrderStatus GROUP BY orderstatusdate) os
      ON m.moid = os.moid
JOIN OrderStatusType ost ON os.ostid = ostid.ostid
WHERE m.ordertype in (2,3)
```

# 9) Answer the Business Question

- The BI report should answer the original business question
  - The <u>What</u> in Step 1

- This report is shared with the stakeholders who asked the original business question
  - The <u>Who</u> in Step 1

- If they accept the BI report as an answer to their question, then this is ready to move to production. → SUCCESS
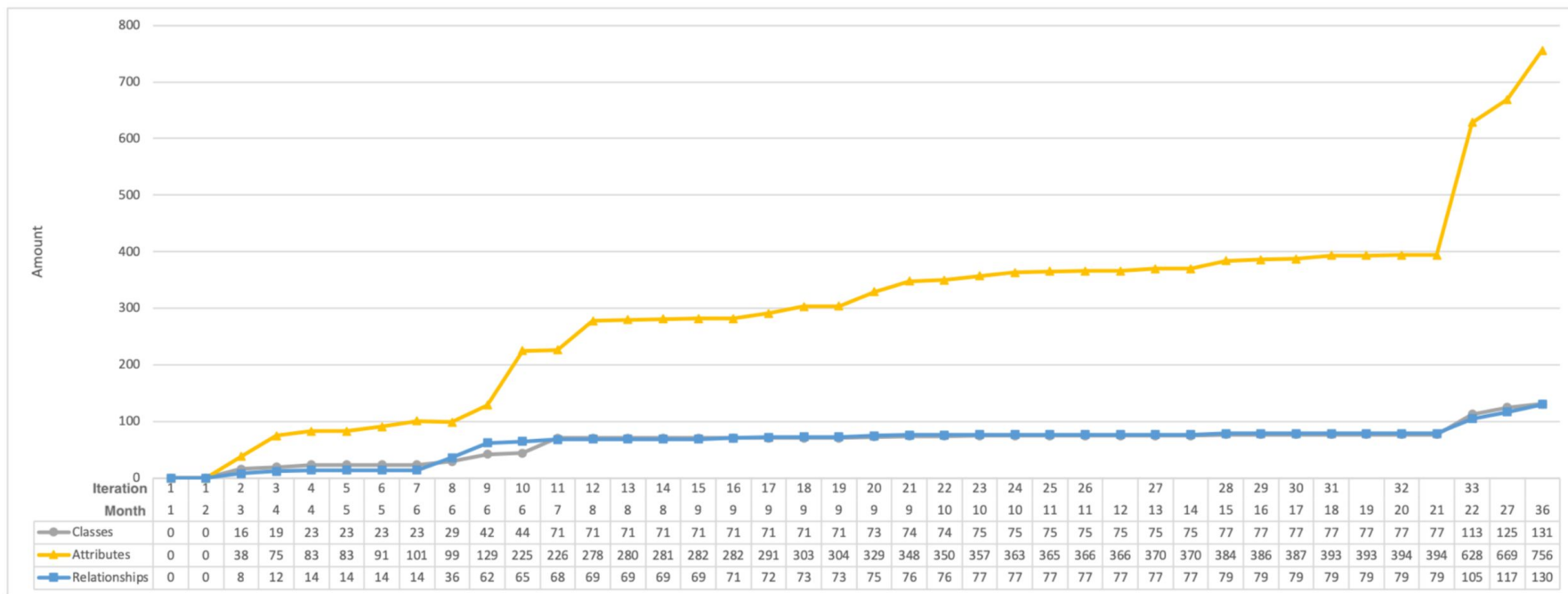
data.world

# 10) Move to Production

- Is the data live (virtual) or ETLed (materialized)?

- Common refresh schedules are daily, weekly, monthly or on demand.

- Determine extract time window:
  - Is the cache going to update the entire extract?
  - or is only yesterday's data going to be pushed to the cache?
  - or last weeks?

data.world

# The E-Commerce Case Study Results

- Goal of First iteration: replicate most trusted BI report, the daily sales report that all C-level executives viewed every morning.

- 3 Business Users, 2 IT users and 1 Knowledge Scientist involved in 2 months.

- Current daily sales report was being generated by one SQL query that a business user would execute every morning.

- The Knowledge Capture phase revealed that the daily sales report encompassed 16 concepts, 38 attributes and 8 relationships. The customer was surprised to see how much knowledge was "hidden" within just one report.

data.world

- **Rapid Growth Phase (Month 1 - 7)**

- **Consistent Growth Phase (Month 8-22)**

- **Independent Growth Phase (Month 23-present)**

# How can we best combine people and technology to improve data integration?

Success depends on this role

Snowball Effect

qual Process

Maintenance & Evolution

Ontology Expressivity

*Interested? Hiring! Looking for research partners!*

juan@data.world
@juansequeda

Ontology and Mapping tools designed for non-semantic aware users

**THANK YOU**

Tools should be designed in conjunction with a methodology

What can be (semi)-automated?

**Knowledge Scientist (People)**

**Pay-as-you-go Methodology (Process)**

**Human vs Machine in the Loop**

**Hybrid Data Cloud (Tools)**

**Dataset Search**

**Nulls... so what?**

**Data Quality & Mappings**

**Learn Mappings from Source Queries**

**SQL Query Log Analysis**

data.world