

# Beauty & the Beast - FAIR Linked Data and the Reality (a pharma perspective)

*Martin Romacker, Data and Information Architect  
Scientific Solution Engineering and Architecture (S2EA)*

*Data & Analytics*

*Roche Innovation Center Basel*

*PLDN 10 Jaar Conference, 27<sup>th</sup> September 2022, Hilversum, Netherlands*



Platform Linked  
Data Nederland

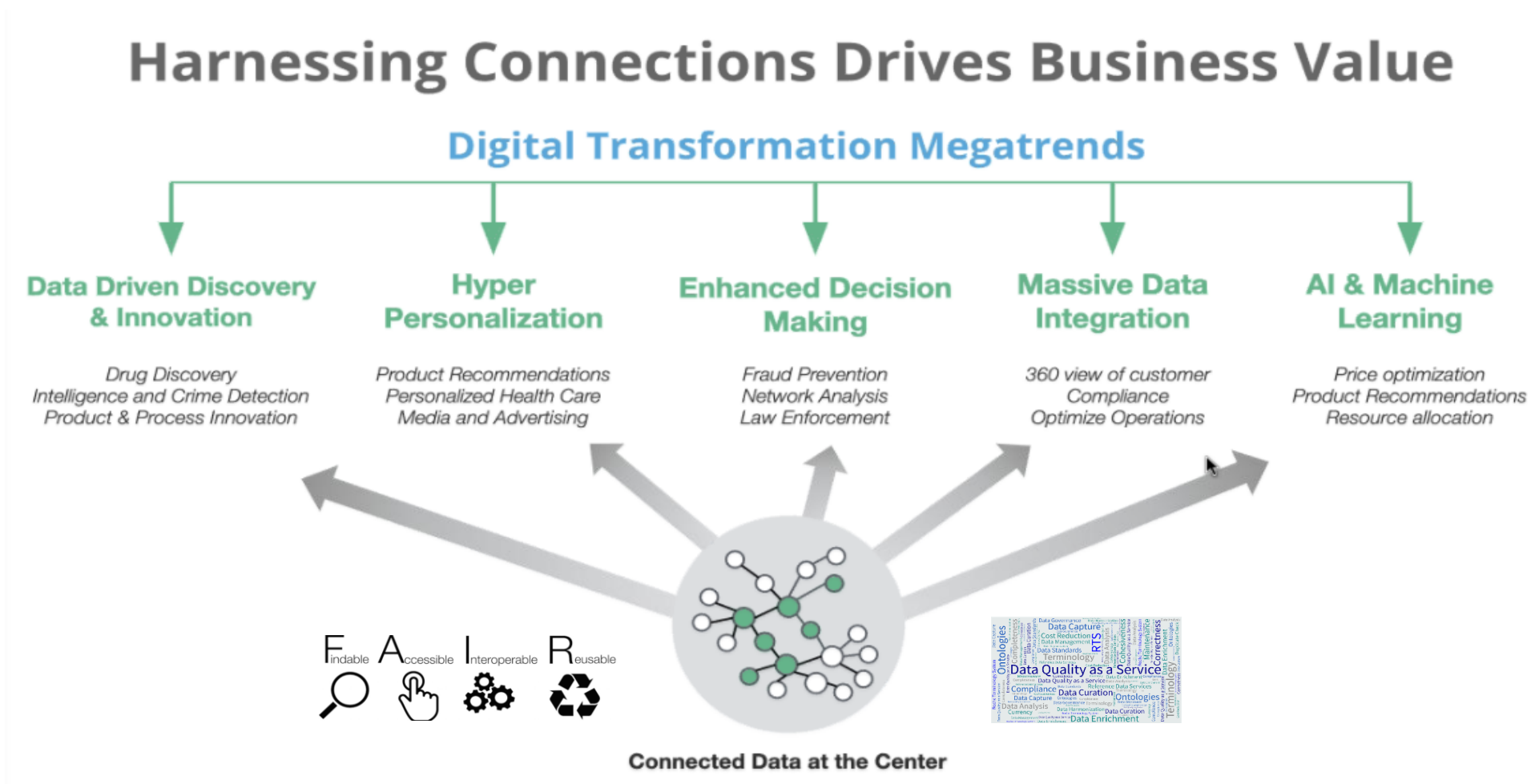
**Roche** *pRED*

# **Digital Transformation & Management of Data Assets**

FAIR plus Q

# Digital Transformation

## Megatrends & Data Management Strategy



The Semantic Web is Dead - Long Live the Semantic Web!

Source: [Rik van Bruggen, Neo4J](#)

# Data as an Asset

## *True Costs of Data Management*

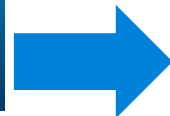


### Planned/ Visible Costs

- FTEs creating Data Asset
- Material procurement (sample, reagent, compounds etc.)
- Infrastructure

### Unplanned/ Invisible Costs

- ETL processes
- Searching & accessing
- Data Cleansing
- Data Curation/ Semantic Data Integration
- IT Infrastructure supporting unplanned activities



Backcharge the costs for processing to the data producers

# **Information Procurement**

## Transformationless Integration of Data Assets

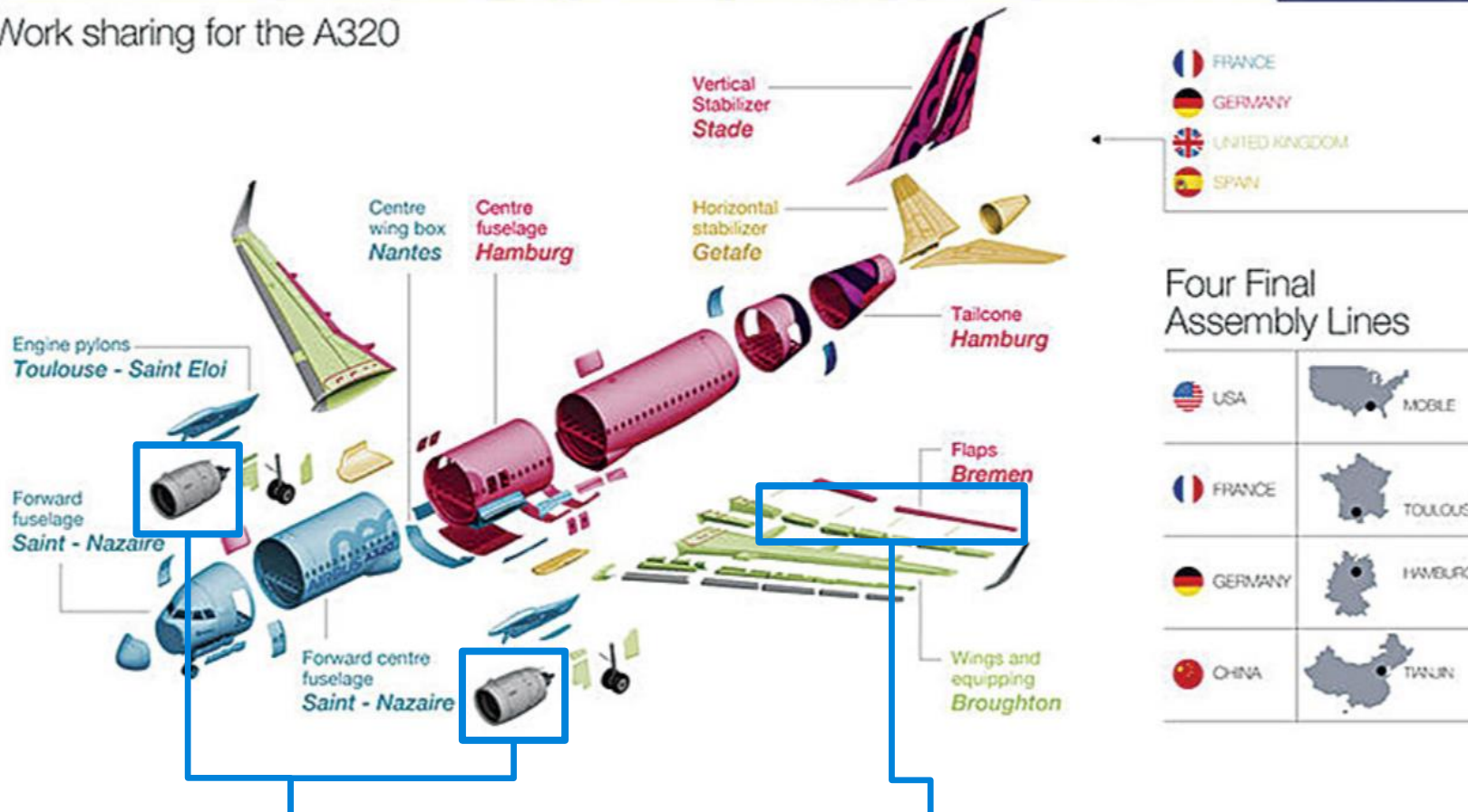


# Manufacturing

## Components & C-Parts Procurement

### Airbus A320, a truly global programme

Work sharing for the A320



- highly complex product
- highly distributed process
- vertical integration (CMO)
- assembly in different global geographic regions

Are data-driven data management/ information assembly processes different?

Contract Manufacturing Organization

Technical Connections/ Interfaces

# Information Procurement

## *Data-Driven and Knowledge-Based Pharma R&D (Ontologies)*

- Information Architecture:  
*Information-centric data organization - semantically sound and meaningful (ontologies)*
- Information Procurement:  
*the effective and efficient process of creating, acquiring and integrating standardized information types into information-driven R&D activities*
  - creation: an in-house activity which will result in a new data asset
  - acquisition: an internalization of a new data asset created by an external organization
  - integration: assembly of internal and external data assets into larger meaningful assets
  - information type: primary building blocks for representation of data assets based on interoperable Minimal Models using community standards (FAIR data principles)

# **Biomedical Ontologies & Terminologies**

## Missing Community Strategy – Intractable Knowledge Space



# FAIRsharing Catalog of Biomedical Resources

## *Proliferation and Fragmentation of Standards*

[Standards](#)
[Databases](#)
[Policies](#)
[Collections](#)
[Add/Claim Content](#)
[Stats](#)
[Log in or Register](#)

### Standards

Contribute by adding a standard | Any problems? Please tell us!

The standards in FAIRsharing are manually curated from a variety of sources, including [BioPortal](#), [MIBBI](#) and the [Equator Network](#).

Manually done-  
no smart interfaces

View as Table | View as Grid

Sort by: Name

Recommended Records

Recommended

Associated Publication?

Showing records 1 - 50 of 1299.

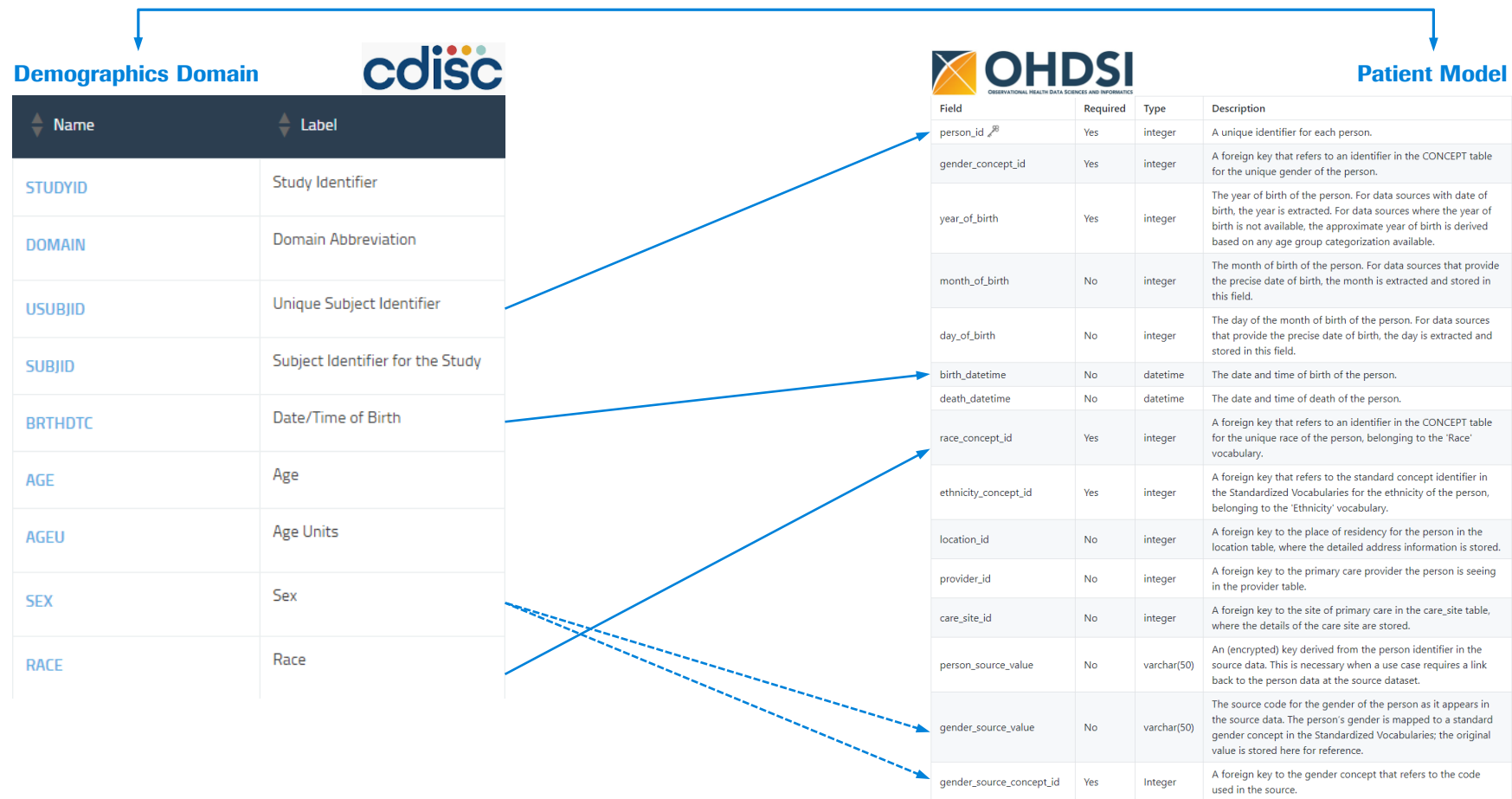
« 1 2 3 4 5 6 7 8 9 10 19 20 21 22 23 24 25 26 »

Registry	Name	Abbreviation	Type	Subject	Related Database	Related Standard	Related Policy	In Collection/Recommendation	Status	
	ABA Adult Mouse Brain	ABA	Standard	None	None	None	None	None		
	Access to Biological Collection Data	ABCD	Standard	<div>Biodiversity</div> <div>Biology</div> <div>Life Sciences</div>	None	<div>GBIF</div> <div>Atlas of Living Australia</div> <div>IPT - GBIF Australia</div>	<div>ABCD</div> <div>EFG</div> <div>ABCDNA</div>	None	<div>TDWG Biodiversity Information Standards</div>	

1299 entries for Standards

# Data Standards & Interoperability Challenges

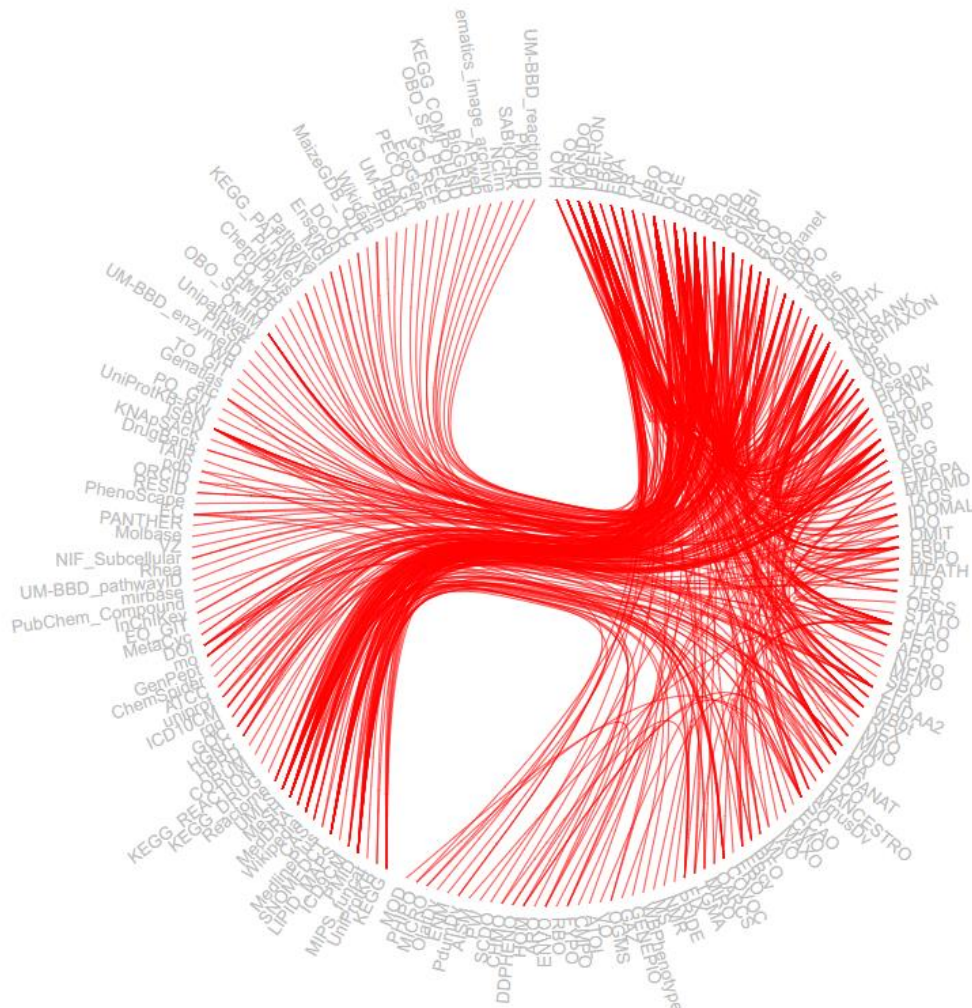
## CDISC vs OMOP/ OHDSI



**Creation of insights & analytics blocked: different model, variables and values**

# EMBL-EBI Ontology Xref Service

## *Creating referential identity by ontology mapping*



Welcome to the EMBL-EBI Ontology Xref Service (OxO).

OxO is a service for finding mappings (or cross-references) between terms from ontologies, vocabularies and coding standards. OxO imports mappings from a variety of sources including the [Ontology Lookup Service](#) and a subset of mappings provided by the [UMLS](#). We're still developing the service so please [get in touch](#) if you have any feedback.

1. Allocating significant resources to inflate a problem
2. Allocating significant resources to reduce a problem (loss of information & interoperability)

# Interoperability for Ontology Mappings

## *RDF standard for a FAIR representation of OM*



### A Simple Standard for Sharing Ontology Mappings (SSSOM)

About SSSOM, A Simple Standard for Sharing Ontological Mappings

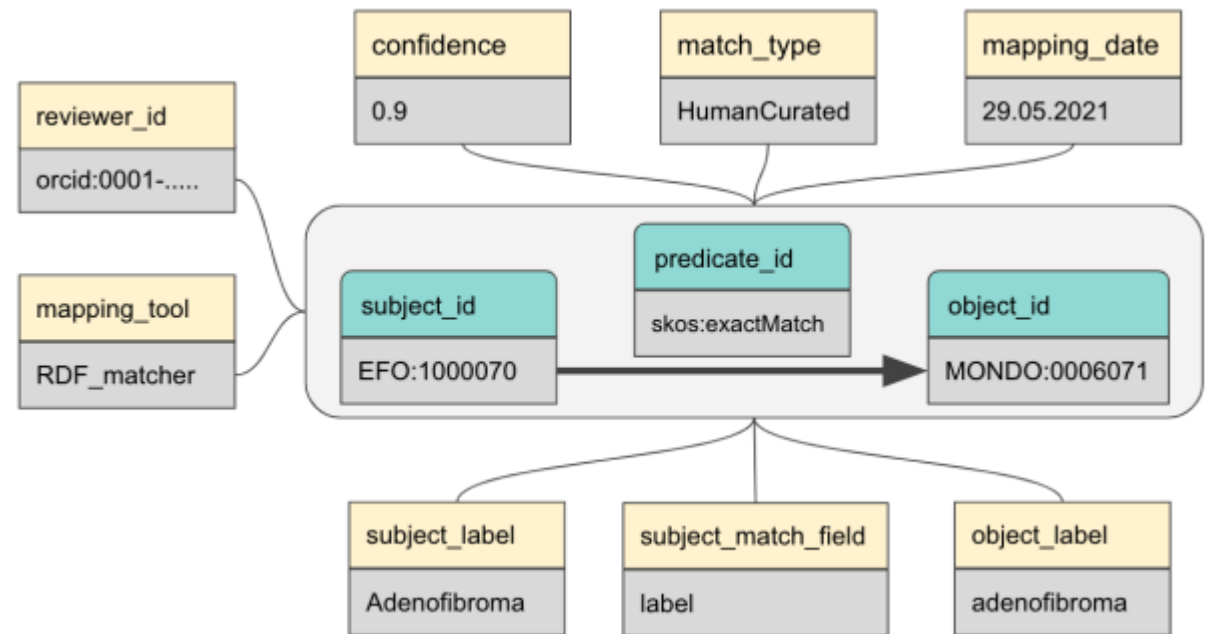
SSSOM is a simple metadata standard for describing semantic mappings:

1. Introducing a machine-readable and extensible vocabulary to describe metadata of mappings.
2. Defining an easy to use table-based format that can be integrated into existing data science pipelines without the need to parse or query ontologies, and that integrates seamlessly with Linked Data standards.
3. Implementing open and community-driven collaborative workflows designed to evolve the standard continuously to address changing requirements and mapping practices.
4. Providing reference tools and software libraries for working with the standard.

A SSSOM mapping comprises three major components:

1. The **mapping** itself, that is, a triple `<subject, predicate, object>` that reflects a correspondence of a `subject` entity, for example a class in an ontology, to an `object` entity, for example an identifier in some database, via a semantic mapping `predicate`, such as `skos:exactMatch`.
2. A **mapping justification**, which the process or activity that led us to consider the mapping to be correct or reasonable (typical examples: labels match exactly; two classes are logically equivalent; a domain expert determined that two terms reflect the same real world concept).
3. **Provenance metadata**, including information about `author` and `mapping_tool`.

[Reference: SSSOM](#)



Not fully FAIR (dct:creator & dct:created)  
No guidelines on property labels

[Linked Open Vocabularies](#)

# **Digital Transformation & FAIRification at Scale**

Industry Approach – Vision & Reality

# **FAIRification at Scale using Biomedical Ontologies**

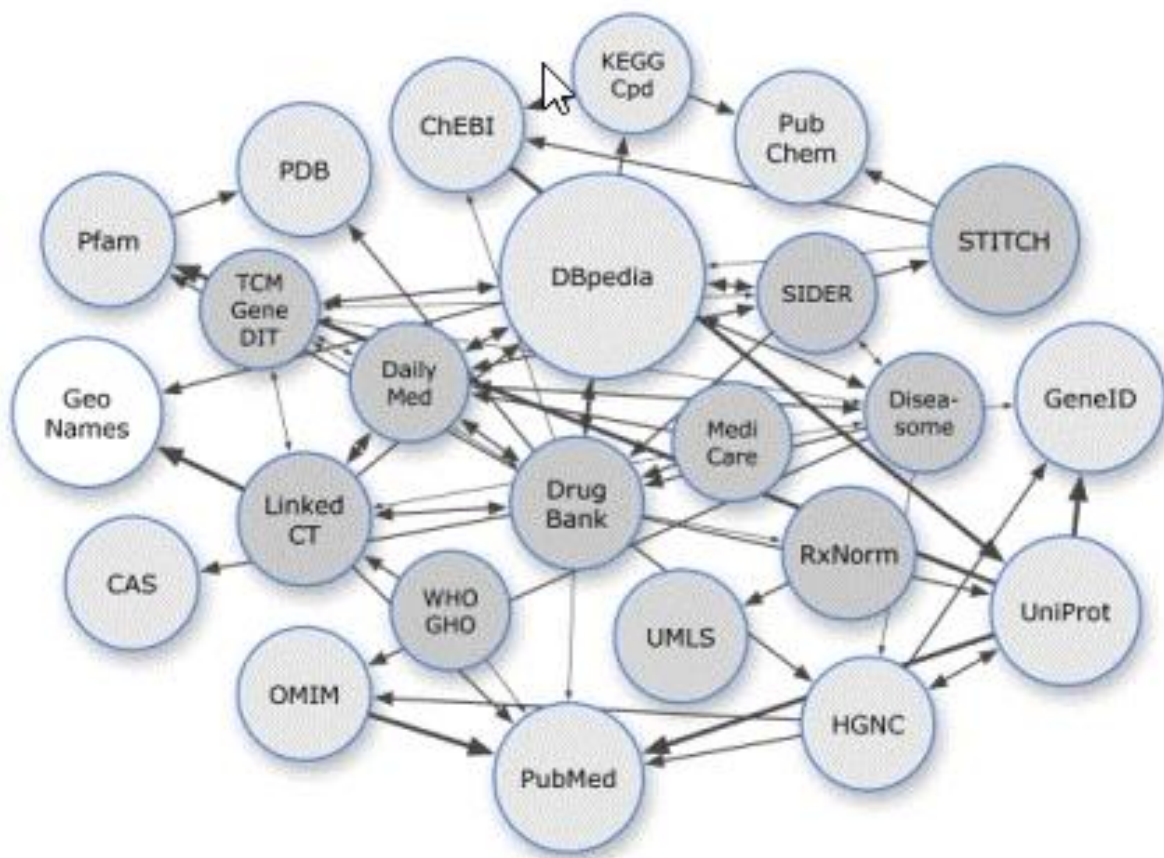
## *Vision*

*An open public-private semantic  
infrastructure of fully standardized  
FAIR applications, services & data*



# Linked Open Data Cloud

## *Linked Open Drug Data Cloud (LODD)*



Data in the Linked Data Cloud is not linked

Linkage comes with referential identity

Referential identity comes with interoperability

LODD is not FAIR

# Digital Data Assets & Data Management

## *The Reality - No Standardization, unlinked, unFAIR*

A	B	C	D	E	F	G	H	I
External_ID	MOLECULE	ORIGIN	Target_class_orig	Target_orig	MoA_orig	MoA_class_orig	Action_o	Pathway_orig
		43 Data Provider 1			Less active enantiomer of naturally-occurring (-)-nicotine	Cholinergic	Agonist	
					is found to be a inhibitor of			
Sata 3	ENVxEa232412	Data Provider 2		PARP	poly (ADP-ribose) synthetase of pancreatic islet cells.			DNA Damage
reposENVxEa3792	ENVxEa37847	Data Provider 3			Hypertension - a beta blocker. A metabolite of			
ZyQ1572291	ENVxEa901222	Data Provider 4	Ion Channel	Ca2+ channel	, ENVxEa34787-xEa0			
Z203042168	ENV0145777	Data Provider 5		MetabotENVpic glutamate receptor				
cENVwdENV07114	ENV07166121	Data Provider 6		mGlu1				
ENV050307	ENV065122	Data Provider 7			Antibiotic betalactam carbapenem			
ENVpr 8949	ENV0232156	Data Provider 8			VEGFR inhibitor. Also inhibits KIT, RET, MET and FLT3			

8 different data providers

different metadata

Free Text

different value domains

Coordination

- Data Curation: creation of evitable additional costs
- Data Curation: creation of evitable project delay
- Data Consolidation/ Data Integration: misallocation of expensive resources (data scientists)
- Information loss (not all information can be recovered)
- Insight creation (data monetization) serverly impacted

# **Why FAIR and Linked Data might fail**

The seven deadly sins

# Seven Deadly Sins of the FAIR Community

## *Deadly Sin #1*

*“The FAIR community fails to make clear what FAIR really is. In particular, the implications of FAIRification on how we work in Knowledge Management and IT projects are not clear.”*

# Seven Deadly Sins of the FAIR Community

## *Deadly Sin #2*

*“There are foundational misunderstandings about the scope of Data FAIRification in particular in the context of Data Quality Frameworks. FAIR Data and High Quality Data are not the same.”*

# Seven Deadly Sins of the FAIR Community

## Deadly Sin #3

*“We cannot really measure FAIR maturity as the FAIR maturity model is almost incomprehensible. This heavily impacts the correct adoption of the FAIR principles”*

The discovery of digital object should be possible from its metadata. For this to happen, the metadata must explicitly contain the identifier for the digital resource it describes, and this should be present in the form of a qualified reference, indicating some manner of "about" relationship, to distinguish this identifier from the numerous others that will be present in the metadata.

In addition, since many digital objects cannot be arbitrarily extended to include references to their metadata, in many cases the only means to discover the metadata related to a digital object will be to search based on the GUID of the digital object itself.



# Seven Deadly Sins of the FAIR Community

## *Deadly Sin #4*

*“The biomedical community does not really converge on standards. We rather increase the chaos instead of harmonizing both in terms of semantic resources and harmonization projects resulting in an intractable knowledge space.”*

# Seven Deadly Sins of the FAIR Community

## *Deadly Sin #5*

*“The FAIR community stays in a bubble. Insiders connect with insiders. The outreach and integration with the business is poor. Besides a vague understanding there is little support for a true break through from management.”*

# Seven Deadly Sins of the FAIR Community

## Deadly Sin #6

*“Scaling up knowledge management based on FAIR resources and standards requires an operational backbone of FAIR data and services – but no community strategy about the basic resources and their maintenance”*

RDF platform

Linked Open Data platform for EBI data

Google

404. That's an error.

The requested URL / was not found on this server. That's all we know.



The Odds of Winning at Poker

LINKEDCT.ORG 10/05/2022 BLOG



When you play Poker, you will probably want to learn the basic rules and odds. There are a number of factors that will influence your game, from your Hand rankings to your Betting options. Read on to learn more. Then, you can start betting! We'll take a look at some of the most important factors. The odds of winning depend on the hand you have. It's important to know these before you place your bet!

<http://www.ebi.ac.uk/rdf/>

<http://lov.okfn.org>

<http://linkedct.org>

# Seven Deadly Sins of the FAIR Community

## Deadly Sin #7

*“Implementation of a landscape of FAIR data, services and application requires a high engagement with key communities such as IT Architects, Master Data Management, Data Managers. Until now FAIR does not speak to them.”*



*Disclaimer: This is an arbitrary listing of tools and vendors without implying any recommendation or preference*

# Seven Deadly Sins of the FAIR Community

*Deadly Sin #X*

*“Using RDF and OWL to build Ontologies as well as creating Knowledge Graphs does not at all prevent you from establishing new data silos or producing unFAIR data.”*

# **Digital Transformation & FAIRification at Scale**

FAIR by Design



# FAIR Playbook for IT Professionals

## *Targeting the Key Enablers*



b|

### FAIR Architecture Playbook

*FAIR by Design*

*A Primer for IT Architects,  
Business Analysts and Software Engineers*



**IT & its professionals as key enabler.**  
**Data Managers should not care about FAIR.**

# FAIR Data & Identifiers

## *Global Unique Persistent Resolvable Identifiers (GUPRI)*



**Globally Unique:** *Uniqueness* means that any identifier refers to exactly one Digital Object. *Global validity* means that every Digital Object should have exactly one identifier for reference where *global* is not limited to our organization but ideally would also include the external universe of discourse.

**Persistent:** An identifier never ever changes. An identifier never gets deleted even if the related Digital Object ceases to exist. The metadata of the identifier should also be maintained.

**Resolvable:** Identifiers are resolved by a service that returns the latest version of the object, including its metadata.

**Opaque GUPRI:** no semantics is encoded in the structure of the GUPRI, and it consists solely of the namespace and an identifier. For example, RTS follows this principle by combining the namespace “<http://ontology.roche.com/>” with a random but unique identifier “ROX1302017050223” to “<http://ontology.roche.com/ROX1302017050223>”. The GUPRI does not reveal any semantically relevant information about the entity it refers to.

Transposing these principles to our organization and establishing FAIR identifier management, we need to define and enforce company-wide or even global policies:

- **Namespace registration:** Provision of a repository and a service supporting the definition and governance of namespaces used for the creation of identifiers.
- **GUPRI policies:** Definition of the format and structure for namespaces and identifiers.
- **Generation/minting of GUPRIs:** Unambiguous creation of unique identifiers by a service.
- **GUPRI resolution service:** Service enabling the resolution of GUPRIs for finding and accessing resources.

### Conclusion:

FAIR applications, services, and data require governance, policies, and infrastructure to manage the identifiers space at the global scale.

**Speaking GUPRI:** There are additional elements in the GUPRI giving the consumer hints about the context of this resource. Table REF offers an example. The namespace “<http://clinical.roche.com/study/>” exposes the semantic type of the resource “Study” in the name. This supports the human readability of GUPRIs. Systems for defining speaking GUPRIs can be very sophisticated<sup>10</sup>.

# **Digital Transformation & FAIRification at Scale**

## Standardization & Capability Stack

# FAIRification at Scale: Capability Stack

## *From Terminologies to Domain Models*

**Terminology Management:** The concepts used in our scientific and technical domains are properly defined, typed and organized in a *Terminology Management System*. Each *concept* is given an unambiguous, complete, *preferred label* and a *textual definition*. The concept is complemented by a rich *synonym set* and *cross-references* linking semantically equivalent concepts in other internal or external repositories.

*Every concept is represented by a global, unique, persistent, and resolvable identifier serving as a reference.*

**Dataset Model Management:** In essence, a *dataset model* describes a fully harmonized representation of a *table-like data structure*. The column headers refer to *metadata elements* (variables, field names, properties, attributes - many different names are used). All the metadata elements are defined in a *Metadata Registry* and share the same rich descriptions as concepts in a terminology management system. The set of all metadata elements forms a *(meta)data dictionary* or a *(meta)data catalog*. When a metadata element is selected as a column header to define a dataset, additional properties are set to determine its *value domain*. Value domains are either *data types* (string, date, boolean, etc.) or terminologies. Value domains establish the constraints for the values occurring in the column of the metadata element.

*Every metadata element is represented by a global, unique, persistent, and resolvable identifier serving as a reference.*

*Every dataset model is represented by a global, unique, persistent, and resolvable identifier serving as a reference.*

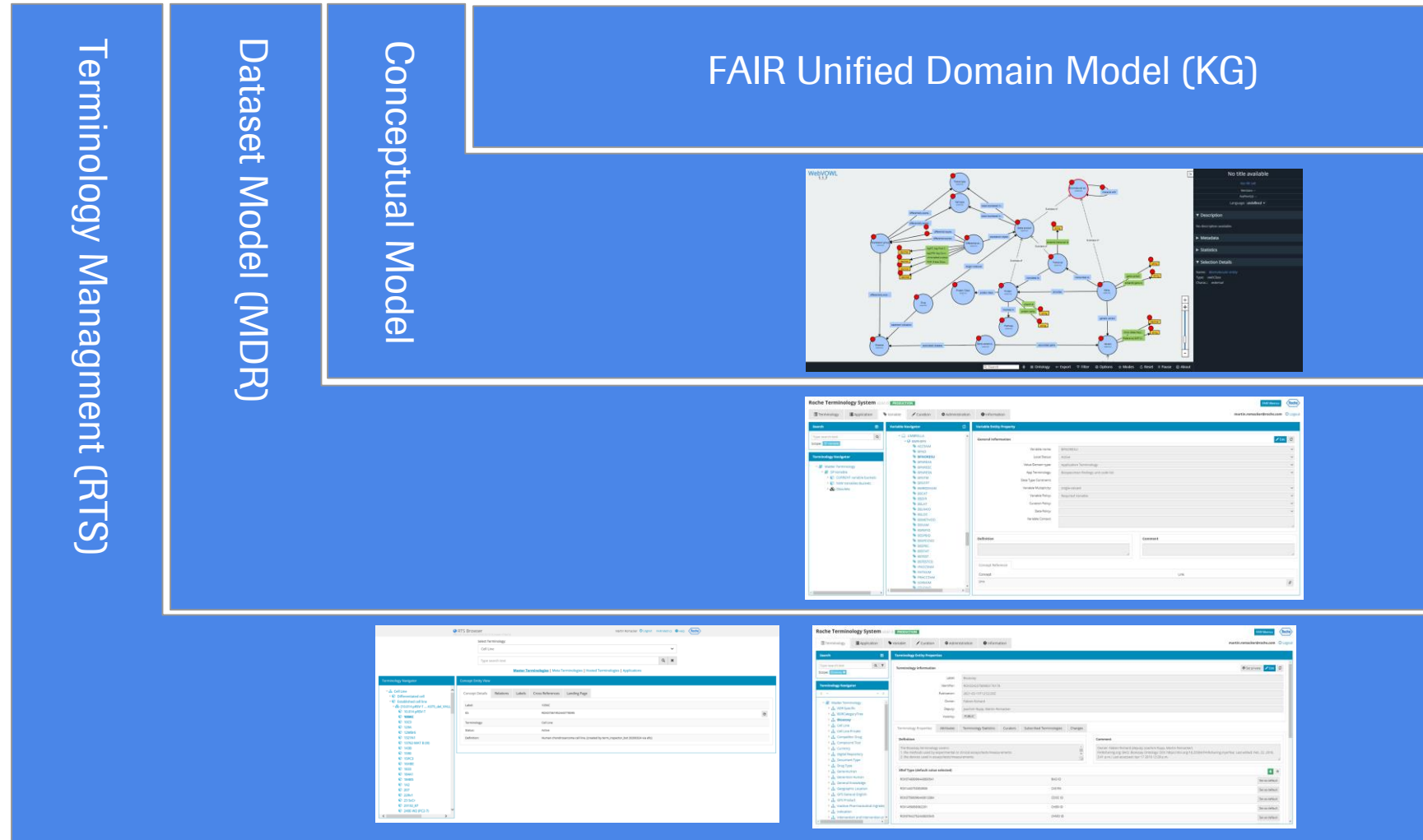
**Conceptual & Logical Model (Domain Ontology) Management:** Following modern data and information architecture approaches, conceptual models support a reasonably grained division of the knowledge space in *data domains* and *subdomains*. In contrast to the table-like dataset models, conceptual models are purpose-driven *Ontologies* representing the *classes* and *properties* of a domain using a directed acyclic graph as a data structure. Domain ontologies can be used as a blueprint for knowledge graphs.

*Every class or property is represented by a global, unique, persistent, and resolvable identifier serving as a reference.*

*Every conceptual or logical model is represented by a global, unique, persistent, and resolvable identifier serving as a reference.*

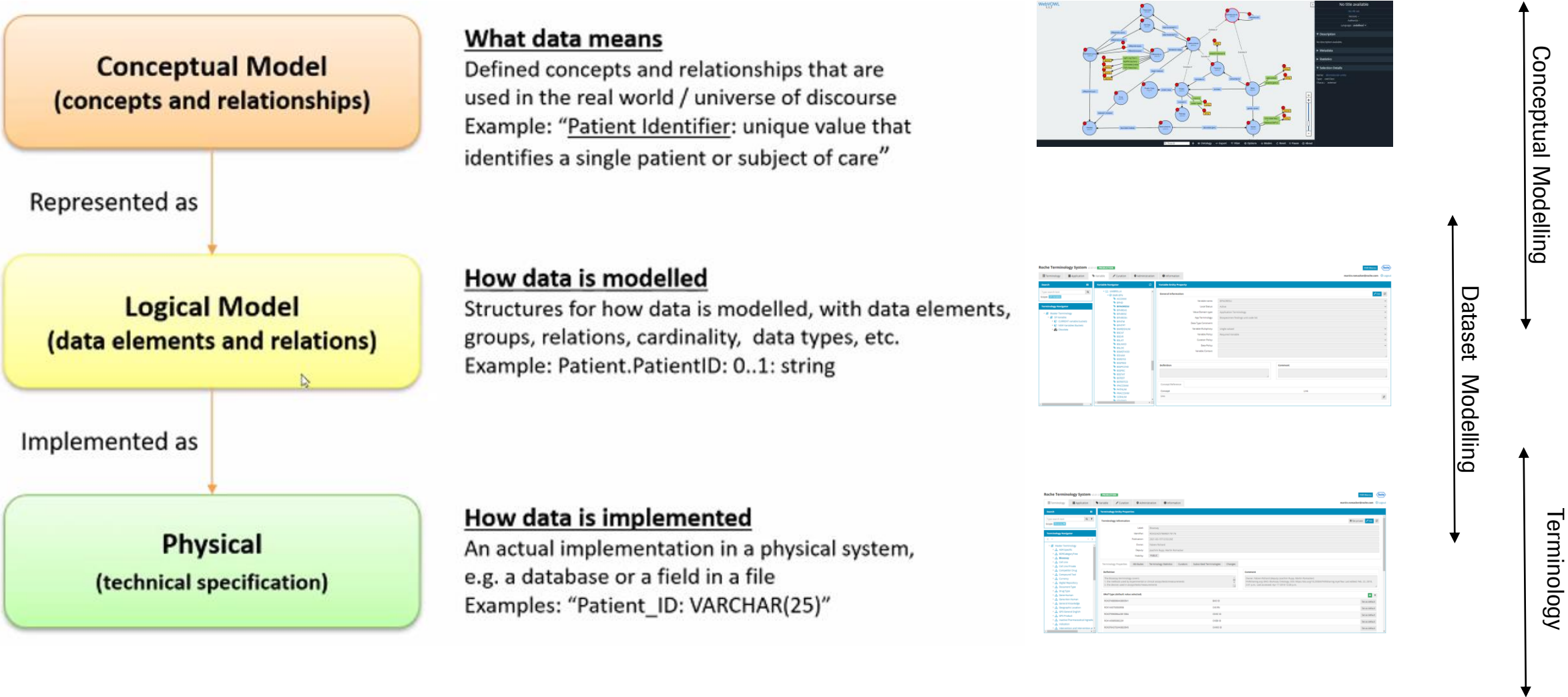
# Semantic Interoperability Hub - Capability Stack

## *Data Management Value Chain - From Terminologies to a Unified Domain Model*



# Data and Information Architecture

## *Mapping Semanti capabilities - Fully FAIR Representation*





# **Semantic Interoperability Hub**

Terminology Management, Metadata, Dataset Model & Ontology

# Reference Data Services for Data Management

## Terminology Management - Contextualize Concepts (FAIR)

Roche Terminology System v2.52.0 [PRD]

FAIR Metrics

Terminology Application Variable Curation Administration Information

-> More than 100 productive applications integrated martin.romacker@roche.com Logout

Search

Non Small Cell Lung Cancer

Scope: Indication

Terminology Navigator

- Basaloid Carcinoma
- Borderline Malignancy Carcinoma
- Breast Carcinoma
- Carcinoma ex Pleomorphic Adenoma
- Carcinoma in Situ
- Carcinoma of Unknown Primary
- Embryonal Carcinoma
- Epithelial-Myoepithelial Carcinoma
- Gastric Pylorus Carcinoma
- Head and Neck Carcinoma
- Hepatocellular Carcinoma
- Large Cell Carcinoma
- Lung Mucoepidermoid Carcinoma
- Lymphoepithelial Carcinoma
- Mammary analogue secretory carcinoma
- Metaplastic Carcinoma
- Metastatic Carcinoma
- Mucinous Carcinoma
- Mucoepidermoid Carcinoma
- Neuroendocrine Carcinoma
- Non Small Cell Carcinoma
  - Non Small Cell Lung Cancer**
    - Adenosquamous Cell Lung Carcinoma
      - ALK mutation positive non small cell lung cancer
      - EGFR mutation positive non small cell lung cancer
      - HER2 mutation positive non small cell lung cancer

Application Navigator

- E-Sample Flow
- e21 HCP Portal
- EpiCX
  - Approved Indication
  - Country
  - CP Indication Mapping
  - Global Product Name
  - Product
  - Scientific Area Indication
  - Trademark Name AU
  - Trademark Name BR
  - Trademark Name CA
  - Trademark Name TW
  - Trademark Name US
    - Actemra
    - Alecensa
    - ANTI-HER2 TDC
    - astegolimab
    - Avastin
      - Cervical Cancer
      - Colorectal Cancer
      - Glioblastoma
      - Malignant Mesothelioma
      - Non-Small Cell Lung Cancer**
      - Ovarian Cancer
      - Renal Cell Carcinoma
    - balovaptan
    - basmisanil
    - Boniva
    - Cadherin-11 mAB
    - Cathflo Activase
    - CEACAM5 CD3 TCB

Concept Entity Properties

General information

Label: Non Small Cell Lung Cancer

Terminology: Indication

Status: Active

Identifier: ROX1305277804386

Definition

A group of at least three distinct histological types of lung cancer, including squamous cell carcinoma, adenocarcinoma, and large cell

Comment

References

Relations

Application References

Landing Page

Lifecycle

Changes

Label	Language	Source	Label Type	Lexical Type
Non Small Cell Lung Cancer	en	Roche	Synonym	prefLabel
Cancer, lung, non small cell	en	PIP	Synonym	altLabel
Cancer, non small cell lung	en	Roche	Synonym	altLabel
Carcinoma, Non Small Cell Lung	en	Roche	Synonym	altLabel
Carcinoma, non small cell lung	en	Roche	Synonym	altLabel
Carcinoma, non small cell lung cancer	en	Roche	Synonym	altLabel
Carcinoma, Non-Small-Cell Lung	en	Roche	Synonym	altLabel
Non small cell lung cancer	en	ADIS, TPP	Synonym	altLabel
Non small cell lung cancer (NSCLC)	en	Roche	AcroDefinition	altLabel

# Reference Data Services for Data Management

## *Metadata Registry/ Dataset Models – Metadata Harmonization (FAIR)*

Roche Terminology System v2.52.0 [PRD]

FAIR Metrics

[martin.romacker@roche.com](mailto:martin.romacker@roche.com) [Logout](#)

Terminology

Application

Variable

Curation

Administration

Information

Search

Country

Scope: SP Variable

Terminology Navigator

- DM variable
  - Actual Arm Code
  - Age
  - Age Units
  - Animal Status
  - Birth Delivery Procedure
  - Country
  - Date and Time of Death
  - Date Time of Birth
  - Date Time of Data Collection
  - Date Time of End of Participation
  - Date Time of First Study Treatment
  - Date Time of Informed Consent
  - Date Time of Last Study Treatment
  - Description of Actual Arm
  - Description of Planned Arm
  - Domain Abbreviation
  - Ethnicity
  - Globally Unique Subject Identifier
  - Investigator Identifier
  - Investigator Name
  - Organism Species Subspecies

Variable Navigator

- HDAP Subject
  - DM Domain
    - Age
    - Age in Days
    - Analysis Age
    - Baseline Body Mass Index (kg per m2)
    - Country
    - Date of Death
    - End Date Time of Treatment
    - End Date of Last Treatment
    - Ethnicity
    - Intent-To-Treat Population Flag
    - Link to Layer 2 dataset
    - Race
    - Safety Population Flag
    - Sex
    - Start Date Time of Treatment
    - Start Date of First Treatment
    - Subject Class Identifier
    - Time from Diagnosis to Rnd (years)
    - Unique Subject Identifier
- HDAP Substance Use
- HDAP Variant
- HDAP Vital Signs
- HDPA Tumor Identification
- HGDI
- HTAg
- I2O Knowledge Base
- IDMP

Application Entity Property

General information

Variable name: Country

Value Domain type: Application Terminology

App Terminology: Country Code (Alpha 3)

Variable Multiplicity: single-valued

Variable Policy: Required Variable

Curation Policy:

Variable Context:

Definition

Country of the investigational site in which the subject participated in the trial (GDSR).

Comment

ISO 3166 format.

Concept Reference

Concept

Country

Link

Data Dictionary

Application

Variable Properties

Variable

# Reference Data Services for Data Management



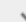

## *Conceptual Model - Purpose-build FAIR Ontologies*

### Model Navigator





- HGDI
- Home Cage Analysis
- I2O Knowledge Base
  - I2O KB Core Model
    - Pathway
    - Genetic variation
      - Reference SNP Cluster Identifier
      - minor allele frequency
    - Gene-disease association
    - Gene variant-disease association
      - associated gene variant
    - Drug
      - target molecule
      - treatment indication
    - Tissue
    - Biomolecule
    - Expression group
    - Cell
    - Disease
    - Gene variant-disease association evidence
      - associated disease
      - associated gene variant
      - p-value
      - odds ratio
      - odds ratio upper 95% confidence interval
      - odds ratio lower 95% confidence interval

### Property Entity View

#### Model Global Properties

Master Concept Identifier:	ROX38009088443943245	
Preferred Label Identifier:	associated disease	
Local Technical Key:		
Preferred Reference URI:		

#### Local Usage Properties

Used at class:	Gene variant-disease association evidence
Target class:	Disease  
Data type:	
Multiplicity:	1..1 

#### Definition

Disease that is part of an association with one or multiple other concepts.

#### Comment

# **FAIR Linked Data**

## Conclusion

# Conclusions

- High-Quality, standardized and linked data: foundation for digitization & insight generation.
- FAIR data principles intrinsically tie Data Management to Semantic Technologies.  
(FAIR data is Linked Data by Design)
- Information Procurement based on FAIR supporting transformationless data integration.
- FAIR is primarily about the \*HOW\* and not only about the \*THAT\* (FAIR maturity indicators).
- Data Management Value Chain: new architectural approaches around data and information.  
Interoperability of terminologies, metadata, dataset models and ontologies is key.
- Data Management Strategy - urgency to build semantic capabilities at community level:  
*open public-private semantic infrastructure of FAIR applications, services and data.*
- It's all about Semantics.

# **FAIR Linked Data**

## Acknowledgements

# Acknowledgements

## *RTS Data Harmonization Service Team*



[Joachim Rupp](#)

RTS Functional Manager, Basel



[Fabien Richard](#)

Terminology Specialist, Basel



[Silvia Jimenez](#)

Terminology Specialist, Basel



[Felix Schwagereit](#)

Scientific Technical Manager,  
Basel



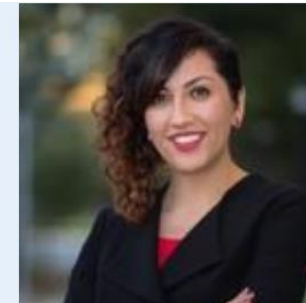
[Pratishtha Duhan](#)

Business Manager, SSF



[Rama Balakrishnan](#)

Biomedical Ontology Specialist,  
SSF



[Shima Dastgheib](#)

Semantic Integrator, SSF



*Doing now what patients need next*