



FAIR Data Implementation Progress Report

Tom Plasterer¹ & Ben Gardner²

¹Strategy, Business Development & Alliances, Oncology
R&D, AstraZeneca, Boston, MA, United States

²Data Standards & Interoperability, Data Office, Data
Science & Artificial Intelligence, R&D, AstraZeneca,
Cambridge, United Kingdom

November 2023



What is Data Centricity?

Data-Centricity puts data at the centre of the enterprise.



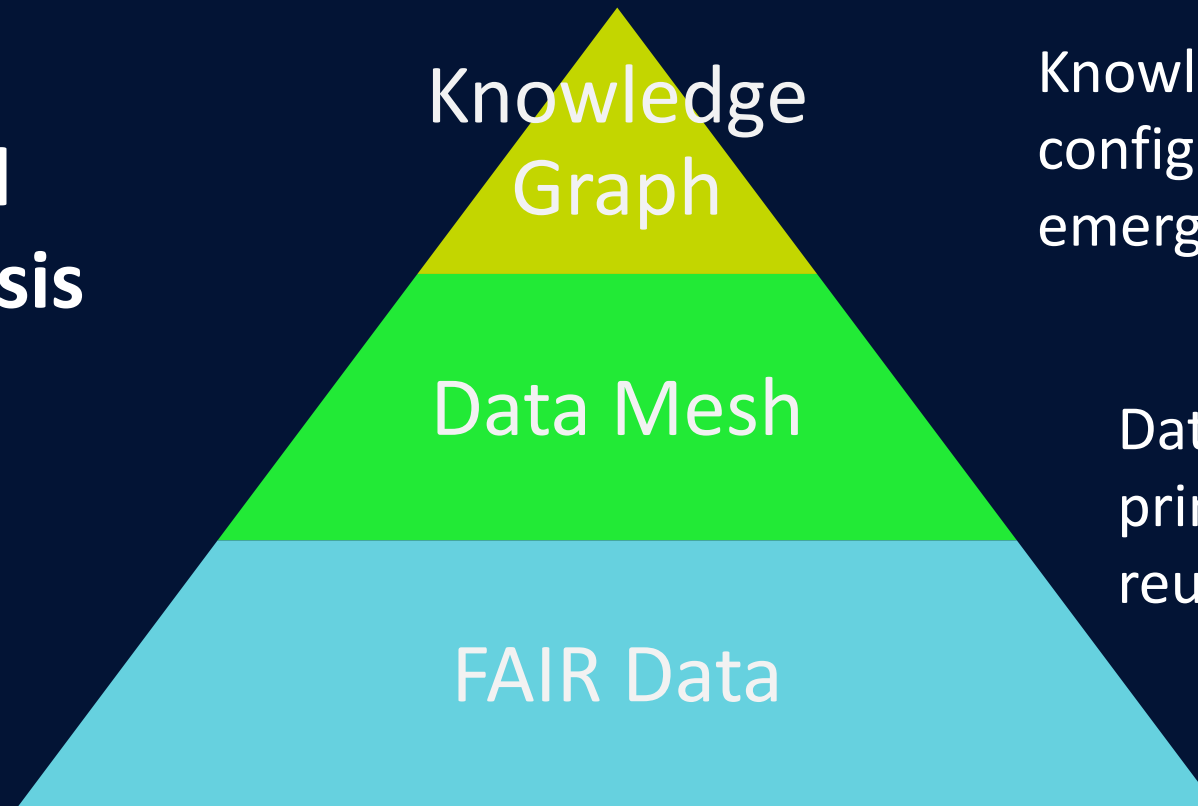
Applications are optional visitors to the data. ([Data-centric manifesto](#))

Data-centricity involves structuring our **data around the science** that we do **rather than the systems** that we use. It promotes data reusability over system-centric design.



Data Centricity Enables Knowledge Applications

FAIR Data is the Foundation for all inquiry and analysis



Knowledge Graph(s) and configurable apps/analytics emerge

Data Mesh utilizes FAIR principles for interoperable, reusable Data Products

Data-centricity can be realized with a commitment to a **FAIR Data foundation** and a loosely-coupled **Data Mesh** manifested in a **Knowledge Graph**.



The value of data centricity

- Random stuff by car
- Ordered by seller
- **Raw data catalogued** and accessible with governance assigned data roles, no insights about what other sellers have
- **Rapid change of quality** and first edition vs expert subject matter knowledge to use

L2



L1



- **Data is catalogued in situ**
- **Randomly distributed** stuff, external data lake
- **Requires expertise** you are looking for
- **Subject matter knowledge** where stuff might be at best access and use



L3

- **Grouping by data category** and spices as a petri dish data lake
- **Improved governance system by**
- **System basis** paid the **provenance** of the **product** data model
- **Control access** at the **enterprise level**
- **Requires average level technical and subject matter knowledge** to use: Data Scientist.

L4



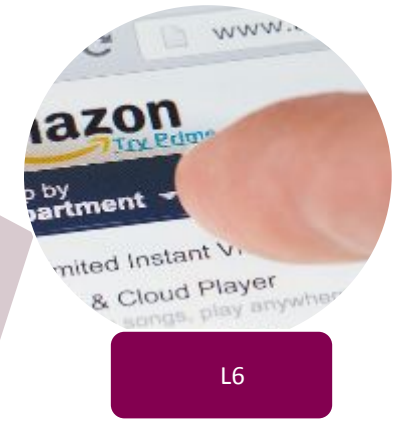
- **Data conformed, integrated, processed and audited** to support specific consumption patterns
- Data is conformed using a **domain level data model**
- **Well organised** Lots of categories
- Data embeds **local master and reference data**
- **Well governed** at the **data level**
- **Well labelled items** including **provenance** (Made where)
- **Creation of analytics ready 'Marts'** enabling self service analytics: Citizen Data Scientist
- **Limited produce**

L5



- **Scale and organisation** at the
- **Data is described in a cross**
- **Everything in one place model**
- **Stores integrated** in a Data Mesh
- **Optimised for enterprise master**
- **societies/communities**
- **Specialist by product category**
- **Data driven display/groupings** -
- **URLs and URLs** summer vs winter; implemented grouping patterns,
- **Context sensitivity**
- **Analytics enabled:**
- **Clear choice of services** all
- **Price comparison** with other retailers

L6

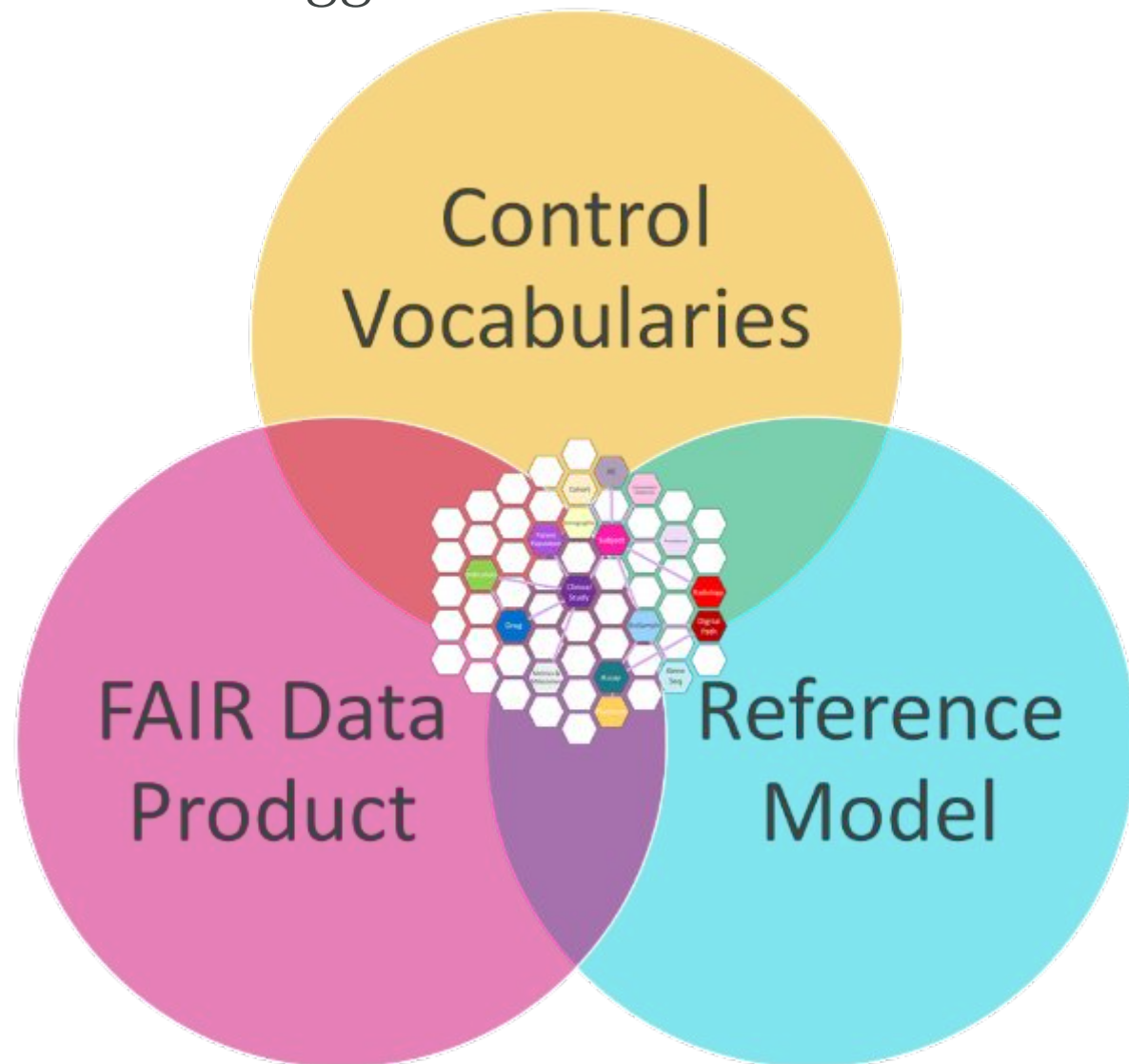


- **Extended digital footprint** no longer physically
- **Data is fully described** using a **knowledge engine** - **ontology** categories/products
- **Data aggregates, business** **correlation** users can **pull** data from **provides information**
- **Automated AI enabled:** AI and **Scale of products and variety** directly **manufactures** without **need for search** helps me find
- **Knowledge enabled** **code** for - I'll know it when I find it searches
- **Amazon subscription services** - schedule delivery, Dash buttons
- **Alexa** - AI guided
- **Services** - Music and video
- **Market place for other vendors**



Building Knowledge Graphs

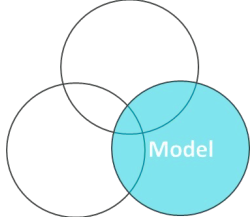
A three-legged stool



“Start with meaning”

Dave McComb, Semantic Arts

Ontology Architecture overview



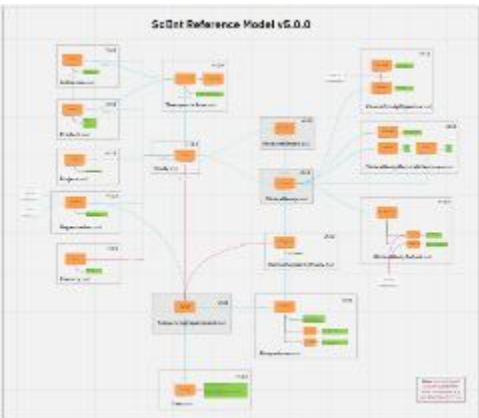
Conceptual Model



Entity Ontology



Application Ontology



Honeycomb



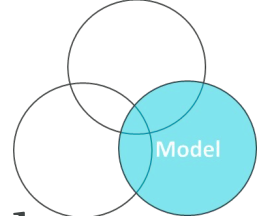
A conceptual model that provides the minimum linking of shared entities used across domain and application models.

Entity Ontologies describe individual concepts and act as building blocks for importing into application ontologies.

Application Ontologies are built to support a specific application.

Honeycombs used to communicate concept of knowledge map, illustrate use case coverage and organic evolution.

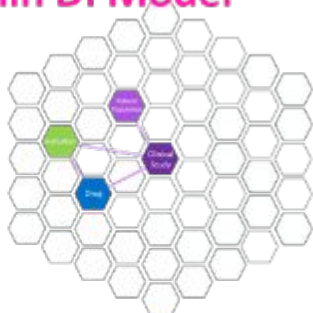




Strategy to organically grow a knowledge map

Increasing the breadth, depth and complexity of questions enabled

Merlin DI Model



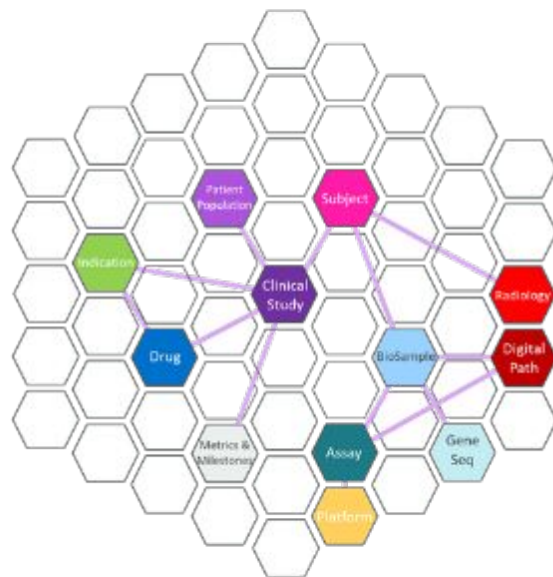
DF+I (PMB) Model



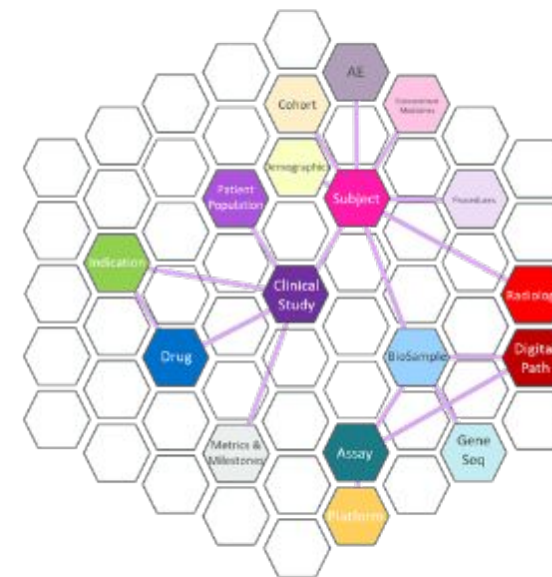
THRIVE model



Combined model



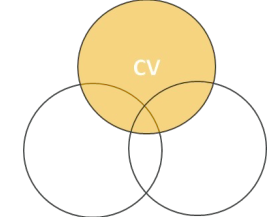
R&D Knowledge Map



Model evolves based on emerging new scientific use cases



Foundational CV's deliver *Quality all the way up*



SKOS-XL

<https://pid.astrazeneca.net/ref/cvname/{ID}>



Collections

A collection of terms derived from existing atomic controlled vocabularies that meet specific application needs/use cases



Foundational CV

Well designed atomic controlled vocabulary built to a common standard

Broad coverage for multiple domains/applications

Decided by SMEs, enabled by specialist curators



Local CV

Well designed atomic controlled vocabulary built to a common standard

Narrow coverage for a domain/application

Decided by use case specific governance, enabled by editorial capability & appropriate tools



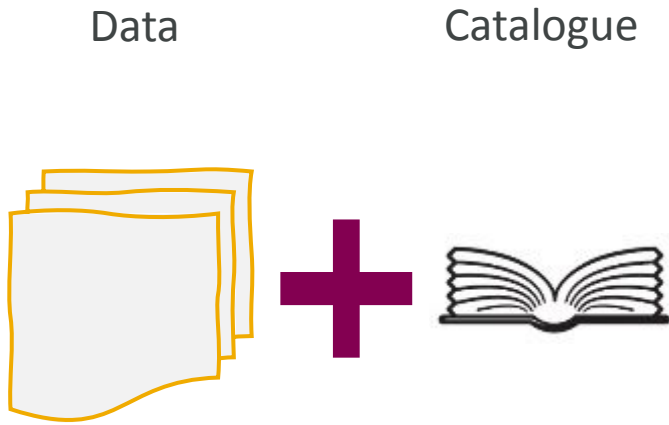
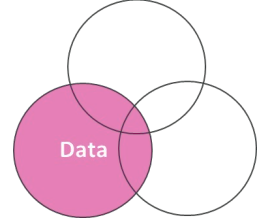
Silo'ed picklists

Uncoordinated list of strings used by an application

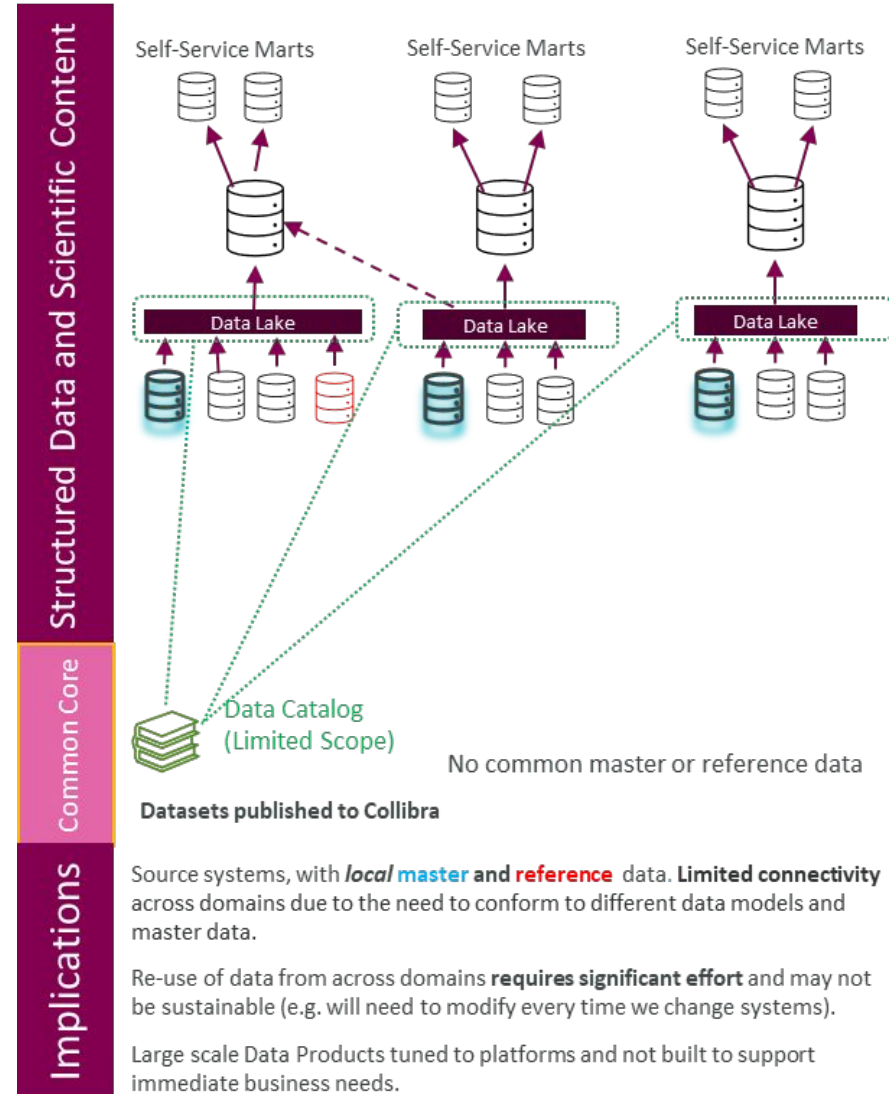


Data

Domain level FAIR (L4) a great first step

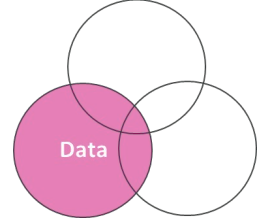


- Find - Data registered in Collibra and tagged with CMM
- Access – Access controlled via Collibra request service
- Reuse – Documentation captured against data

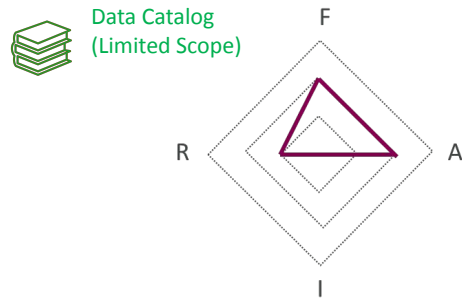


But we are FAR from FAIR

We MUST go further



FAIR metrics (Level 4)



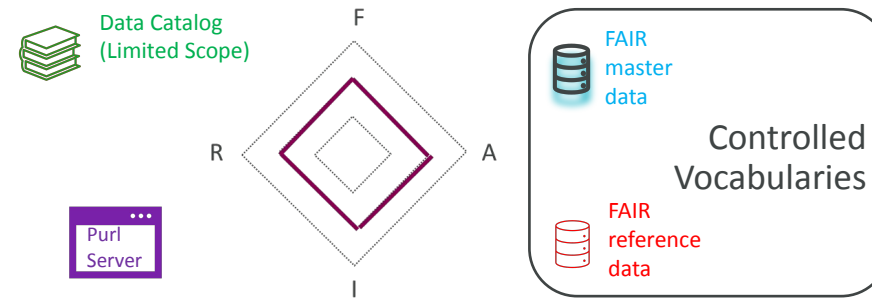
Find - Data registered in Collibra and tagged with CMM

Access – Access controlled via Collibra request service

Reuse – Documentation captured against data

Data record discoverable in Collibra

FAIR metrics (Level 5)



Find - Data registered in Collibra and tagged with CMM

Access – Access controlled via Collibra request service

Interoperable – Data enhanced with shared CV and PIDs

Reuse – Documentation captured against data, includes data dictionary, etc

Data record discoverable in Collibra

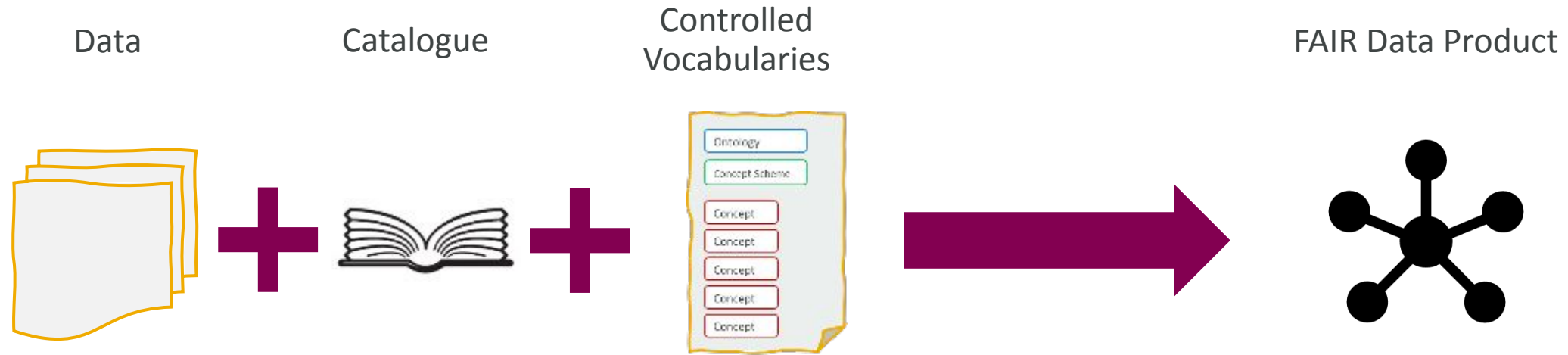
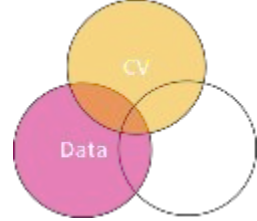
Data enriched to create interoperability

Data is Machine Readable



Enterprise Level FAIR (L5)

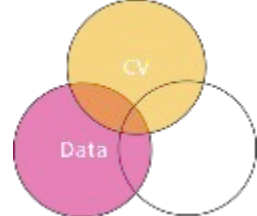
A FAIRe(nough) data product



The minimum viable FAIR Data standard should deliver

Findable	Accessible	Interoperable	Reusable
Registered and discoverable in a Data Catalogue	Mechanism for requesting and receiving data	The data has been aligned to AZ standards where they exist	Documentation describing the constraints associated with using the data
		The PID for each instance in the standard is included to make the data machine readable	Documentation describing the data i.e. data dictionary, schema, etc





Data and Controlled Vocabularies

Putting Interoperability into FAIR

Dirty data

Study	Indication	Drug
D1234C00001	Non small cell lung cancer	Tagrisso
ADORA	NSCLC	Osimertinib
CP11278-CMA33G	Diabetes type 2	Forxiga

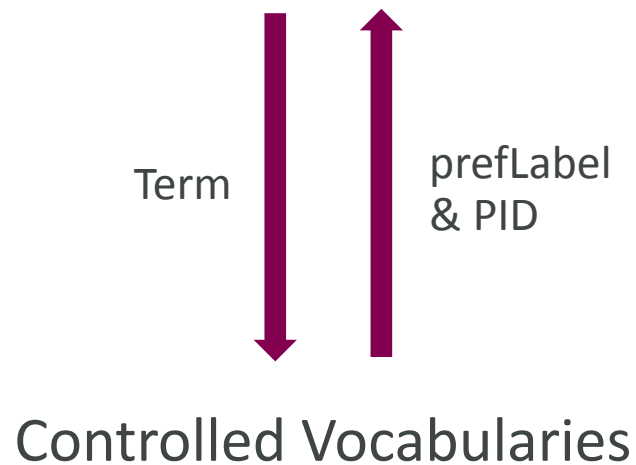
- Inconsistent identifiers & terms
- Column values can be concatenated
- etc



Interoperable data

Study_ID	Study_ID_URI	Indication	Indication_URI	Drug	Drug_URI
D1234C00001	https://pid.astrazeneca.com/1/12345	Non small cell lung cancer	https://pid.astrazeneca.com/Indication/23456	Tagrisso	https://pid.astrazeneca.com/Product/965723
D1234C00012	https://pid.astrazeneca.com/1/48373	Non small cell lung cancer	https://pid.astrazeneca.com/Indication/23456	Tagrisso	https://pid.astrazeneca.com/Product/965723
D4568L00007	https://pid.astrazeneca.com/1/97538	Diabetes type 2	https://pid.astrazeneca.com/Indication/9857	Forxiga	https://pid.astrazeneca.com/Product/853584

↑ prefLabel ↑ PID



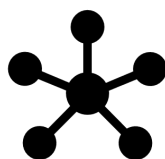
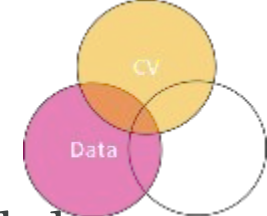
Controlled Vocabularies

- Shared Controlled Vocabularies
 - Enrich with preferred label and PIDs
- Uses common files format CSV, JSON, etc
- Is machine readable, graph enabling and relationally world friendly
- Well documented - Data Dictionary/Data Schema/etc



L5 FAIR Data Products benefit all

Inclusion of PIDs simplifies data integration irrespective of target data model



L5 FAIR Data Product

Study_ID	Study_ID_URI	Drug	Drug_URI
D1234C00001	https://pid.astrazeneca.com/1/12345	Tagrisso	https://pid.astrazeneca.com/Product/965723
D1234C00012	https://pid.astrazeneca.com/1/48373	Tagrisso	https://pid.astrazeneca.com/Product/965723
D4568L00007	https://pid.astrazeneca.com/1/97538	Forxiga	https://pid.astrazeneca.com/Product/853584

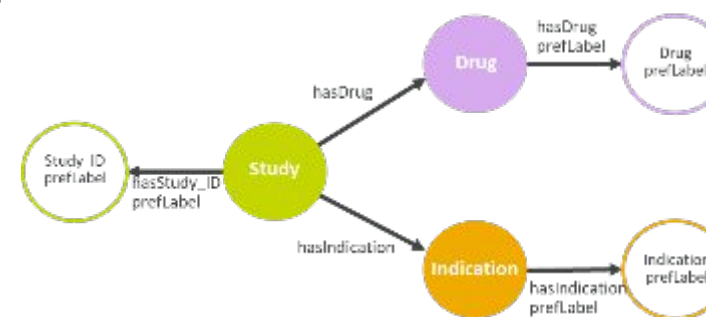


Study_ID	Study_ID_URI	Indication	Indication_URI
D1234C00001	https://pid.astrazeneca.com/1/12345	Non small cell lung cancer	https://pid.astrazeneca.com/Indication/23456
D1234C00012	https://pid.astrazeneca.com/1/48373	Non small cell lung cancer	https://pid.astrazeneca.com/Indication/23456
D4568L00007	https://pid.astrazeneca.com/1/97538	Diabetes type 2	https://pid.astrazeneca.com/Indication/9857

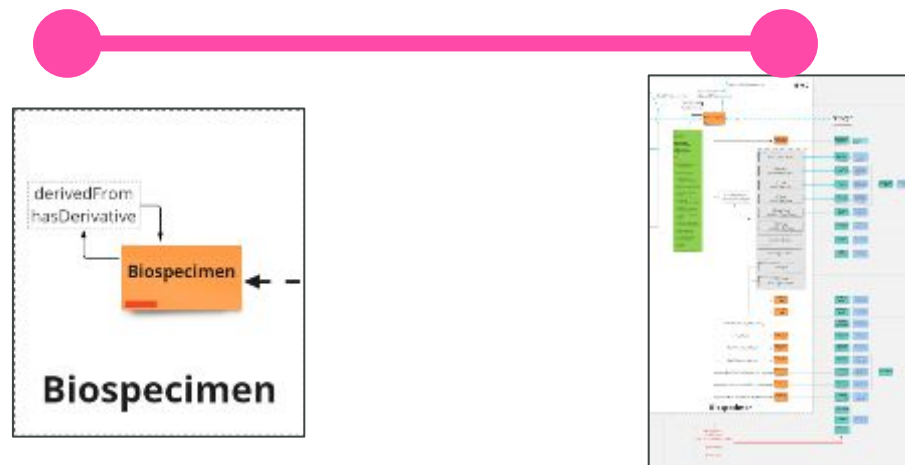
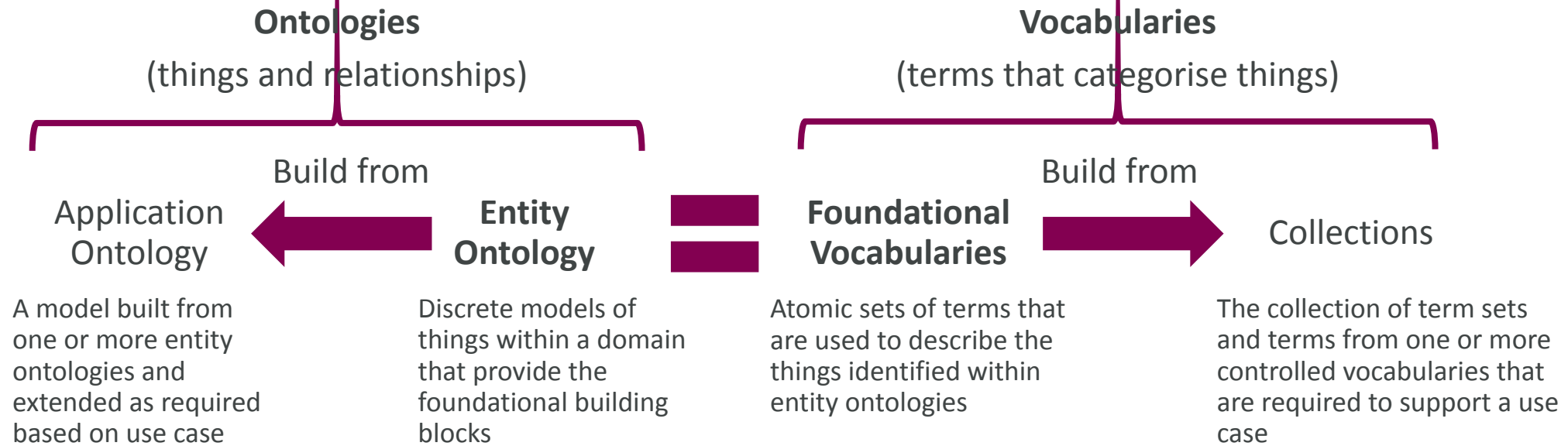
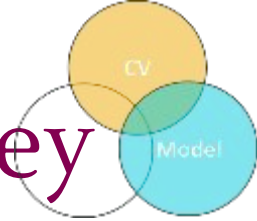
Relational

Study_ID	Study_ID_URI	Indication	Indication_URI	Drug	Drug_URI
D1234C00001	https://pid.astrazeneca.com/1/12345	Non small cell lung cancer	https://pid.astrazeneca.com/Indication/23456	Tagrisso	https://pid.astrazeneca.com/Product/965723
D1234C00012	https://pid.astrazeneca.com/1/48373	Non small cell lung cancer	https://pid.astrazeneca.com/Indication/23456	Tagrisso	https://pid.astrazeneca.com/Product/965723
D4568L00007	https://pid.astrazeneca.com/1/97538	Diabetes type 2	https://pid.astrazeneca.com/Indication/9857	Forxiga	https://pid.astrazeneca.com/Product/853584

Graph

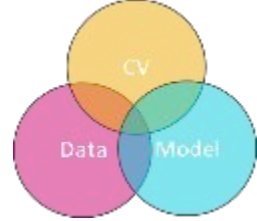


Aligning entities with controlled vocabularies is key



Knowledge Level FAIR (L6)

A FAIR data product minimises the gap between relational and graph worlds

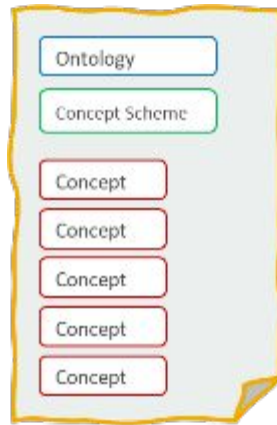


L4 FAIR Data

L4 – Domain level FAIR

- Data conformed, integrated, processed and audited to support specific analytics patterns/enquiries
- Data is conformed using a domain level data model
- Data embeds local master and reference data
- Controls on access at the data level
- Creation of analytics ready 'Marts' enabling self service analytics: Citizen Data Scientist

Controlled Vocabularies

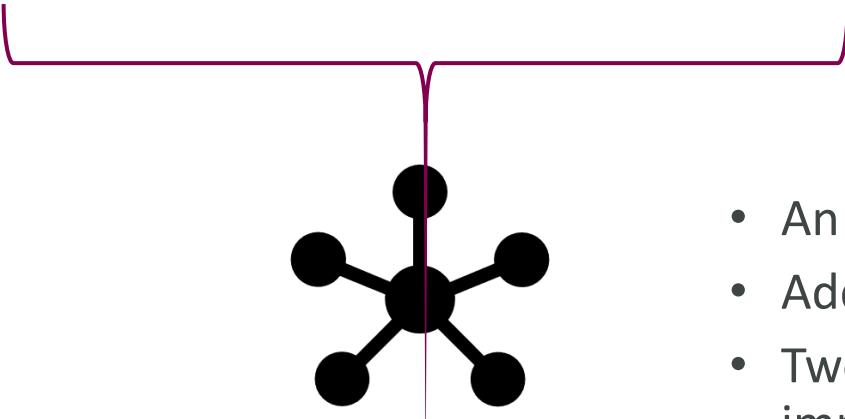


vocab RDF file

Reference Model



Knowledge Map



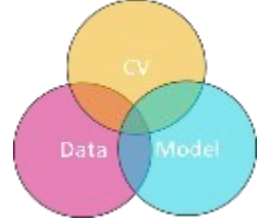
L5 FAIR Data Product

- An enterprise standard
- Adds value for everyone
- Two thirds of the way to graph without impacting non-graph consumers

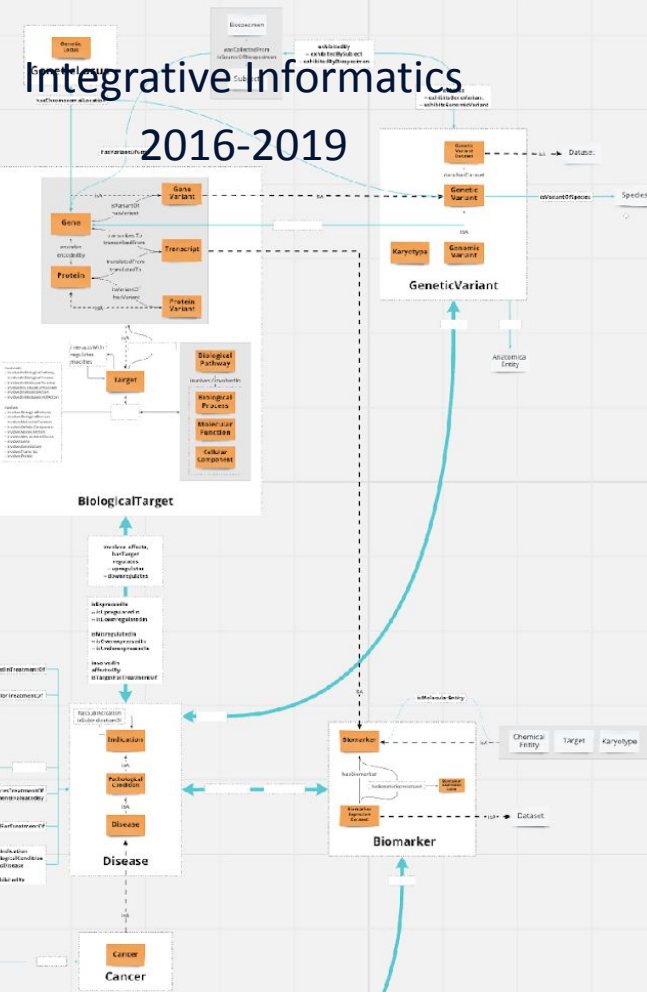
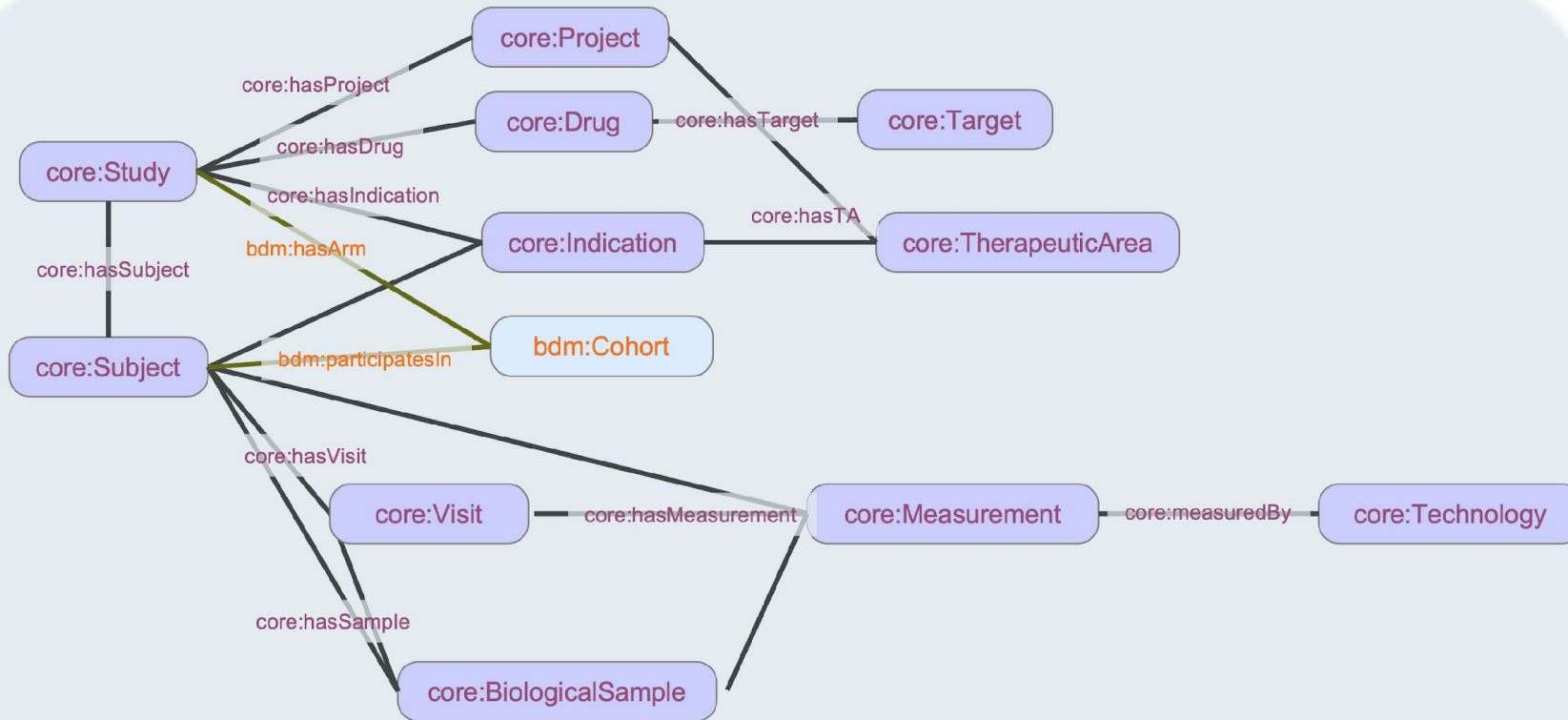


Accepting Change

You may need to give up your tactical solution



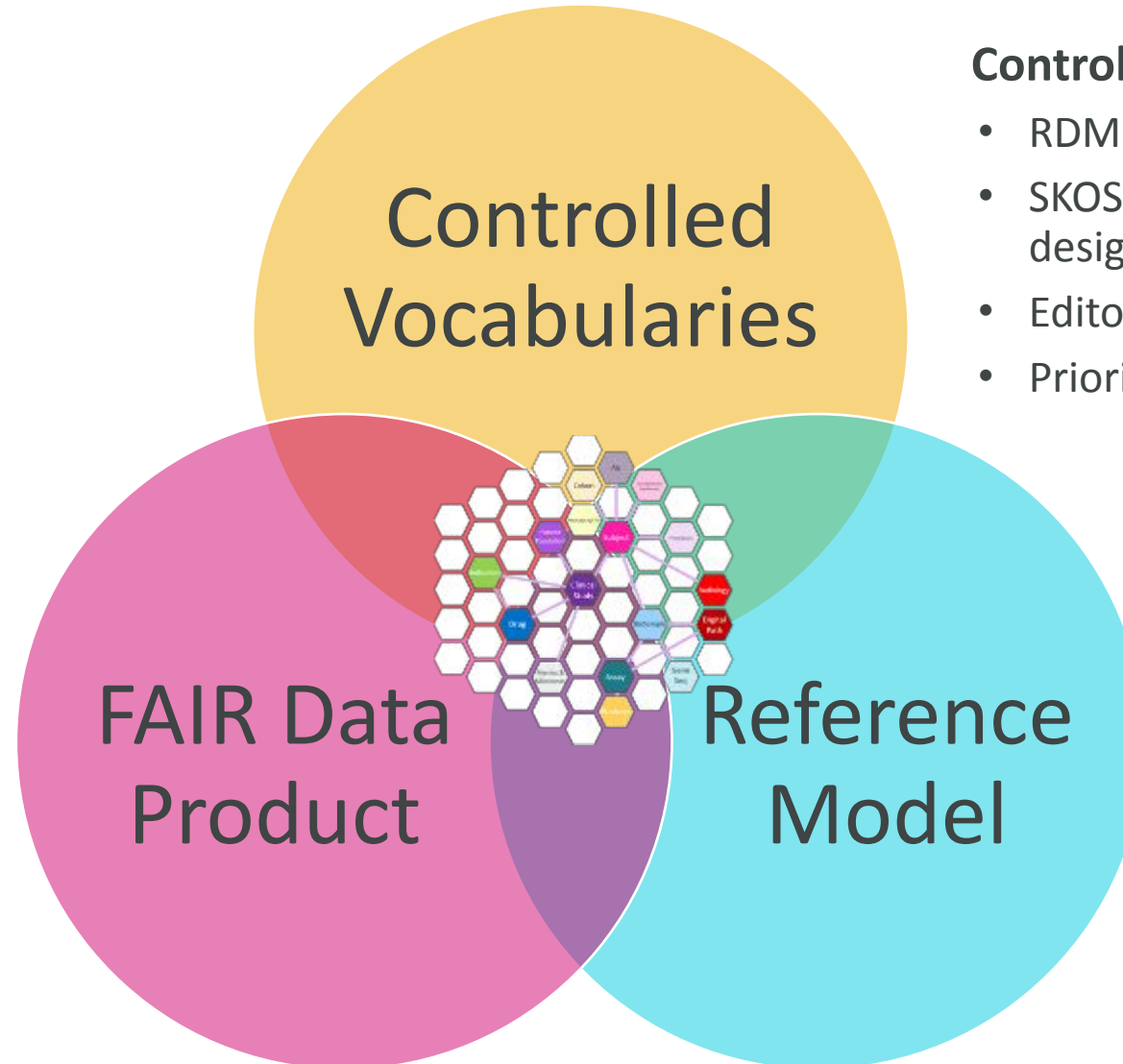
Starting Point: Modeling Business Questions



Bringing it all back together

FAIR Data

- FAIR metrics goal 80% Data FAIR by 2025
- Standards, resources, tooling and governance
- FAIR assessment in place
- PID server and standard patterns agreed



Controlled Vocabularies

- RDM service implemented
- SKOS-XL based CV framework & workflow designed
- Editorial controls in place
- Prioritisation & creation of CV underway

Ontology Architecture

- Scalable and sustainable
- Reusable library of entity ontologies
- Organic evolution of reference ontology



Across AstraZeneca Acknowledgements

Ben Gardner

Mathew Woodwark

Daniel Roythorne

Jon Ison

Nathalie Conte

Nicola Ellingham

Arun Balaji

Induja Mohan

Arinjay Jadeja

Hans Ienasescu

Bhavna Khilnani

Michael Neylon

Rob Hernandez

Derek Scuffell

Varsha Khodiyar

Pablo Porras Millan

John Berrisford

Bijay Jassal

Rafa Jimenez

Philippe Rocca-Serra

Victor Kim

Alex Wood

Linda Zander-Balderud

Antonio Fabregat

Justin Johnson

Mark Reuter

Tom Plasterer

James Holman

Martina Devoti

Stacy Mather

Di Elvers

Colin Wood

Sandra Mc Garry

Gareth Henry

Kerstin Forsberg

Calle Nordmark



 Pistoia Alliance

FAIR Toolkit

1. Metric Tools & Best Practice

2. Training resources

3. Culture change process

4. Use case examples

5. Cost benefit examples

- Adapt for **Life Science industry**
- Leverage **existing** FAIR resources



Questions?



Confidentiality Notice

This file is private and may contain confidential and proprietary information. If you have received this file in error, please notify us and remove it from your system and note that you must not copy, distribute or take any action in reliance on it. Any unauthorized use or disclosure of the contents of this file is not permitted and may be unlawful. AstraZeneca PLC, 1 Francis Crick Avenue, Cambridge Biomedical Campus, Cambridge, CB2 0AA, UK, T: +44(0)203 749 5000, www.astrazeneca.com

