

# Measuring FAIR in AZ: Lessons learned from implementation of FAIR benchmarks

Pablo Porras Millán, PhD

Data Curation Director, AZ

18/04/2023



# F<sub>indable</sub> A<sub>ccessible</sub> I<sub>nteroperable</sub> R<sub>eusable</sub>

Open Access | Published: 15 March 2016

## The FAIR Guiding Principles for scientific data management and stewardship

Mark D. Wilkinson, Michel Dumontier, [...]Barend Mons 

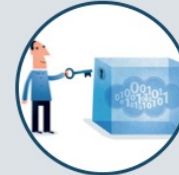
*Scientific Data* **3**, Article number: 160018 (2016) | [Cite this article](#)

**334k** Accesses | **2882** Citations | **1902** Altmetric | [Metrics](#)



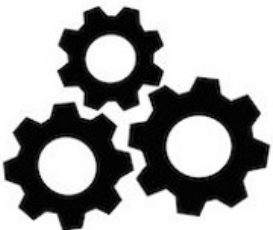
Data and supplementary materials have sufficiently rich metadata and a unique and persistent identifier.

**FINDABLE**



Metadata and data are understandable to humans and machines. Data is deposited in a trusted repository.

**ACCESSIBLE**



Metadata use a formal, accessible, shared, and broadly applicable language for knowledge representation.

**INTEROPERABLE**

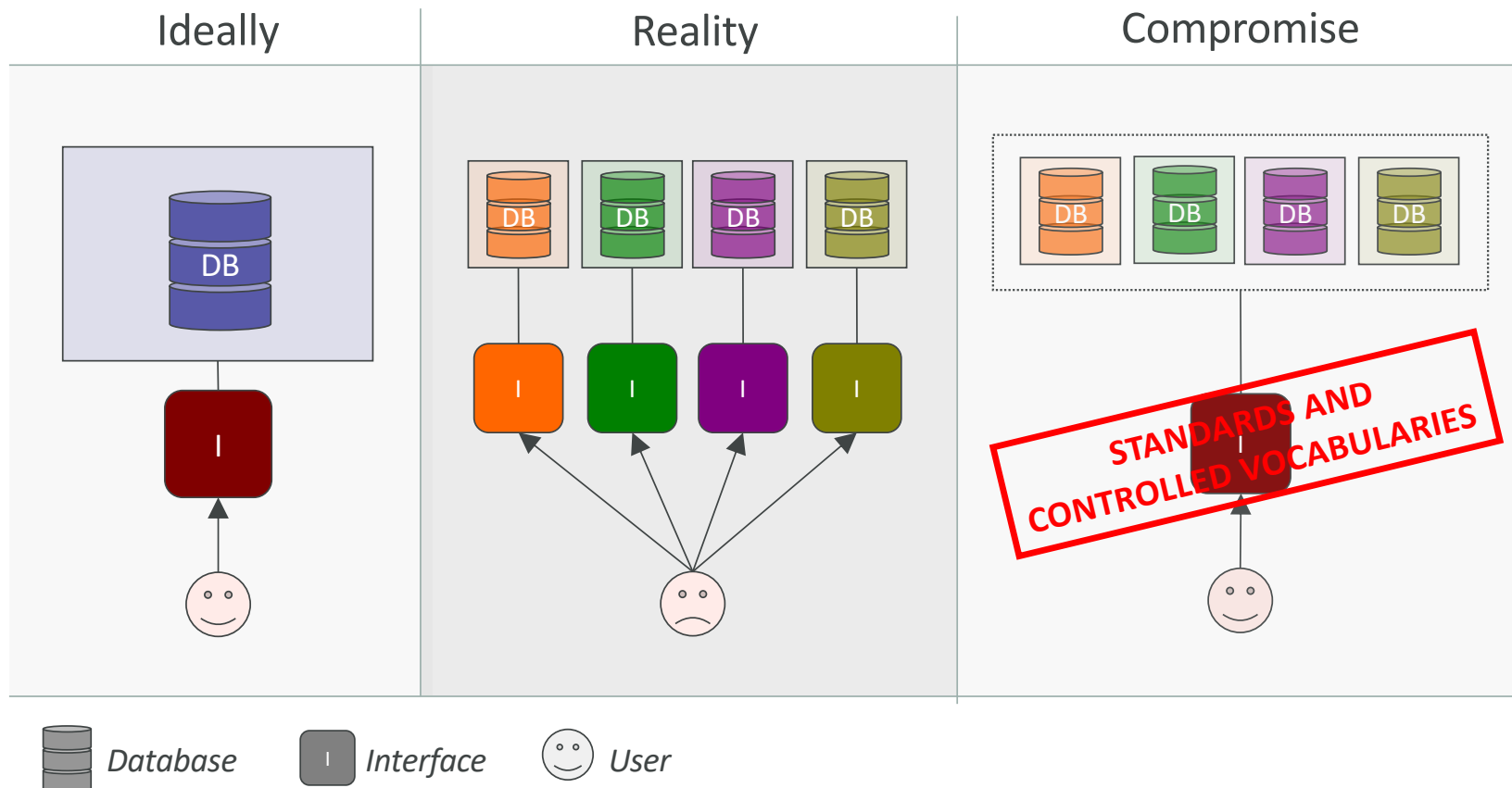
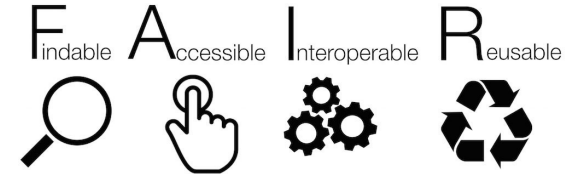


Data and collections have a clear usage licenses and provide accurate information on provenance.

**REUSABLE**

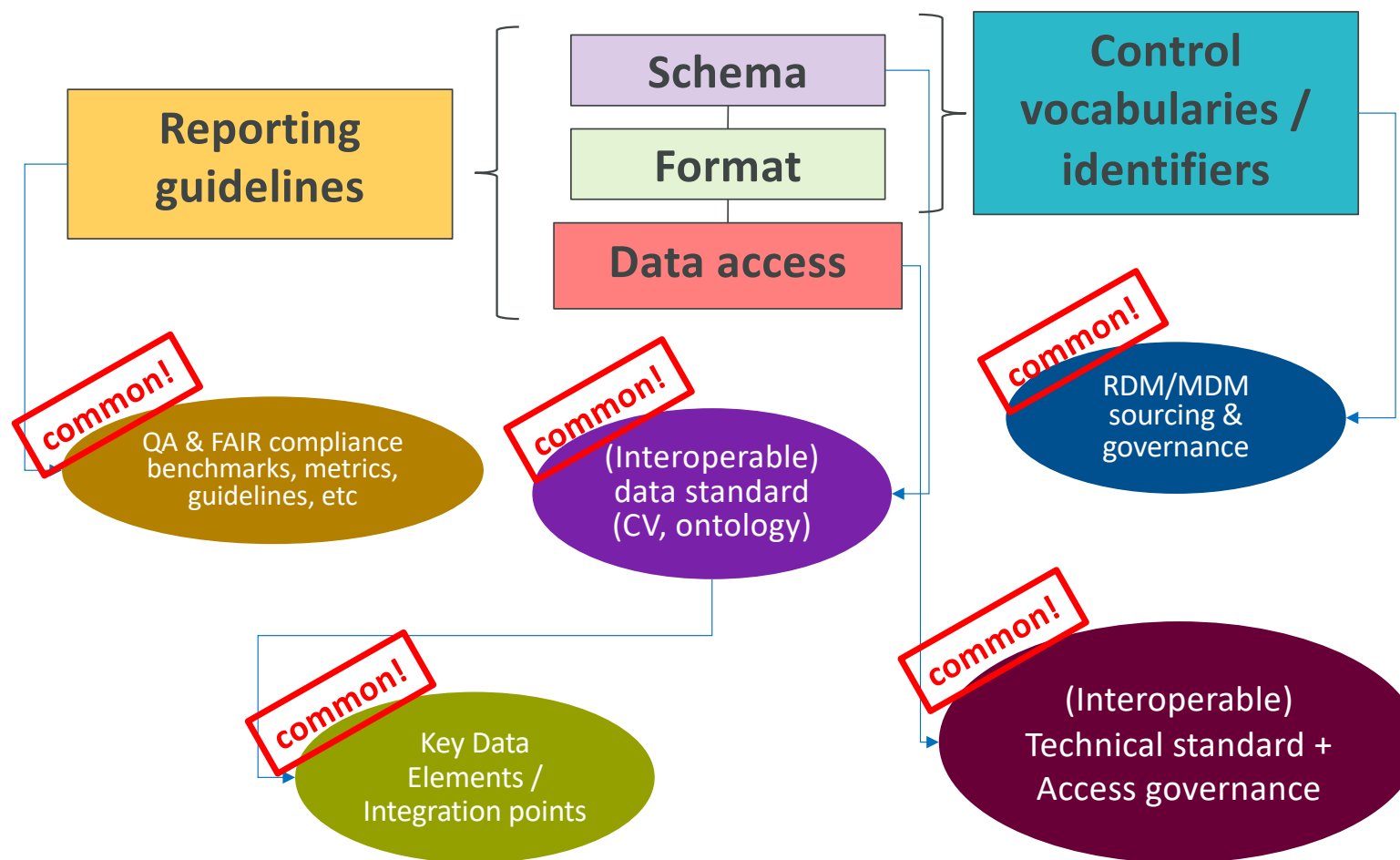


# Application centricity: A problem of data integration

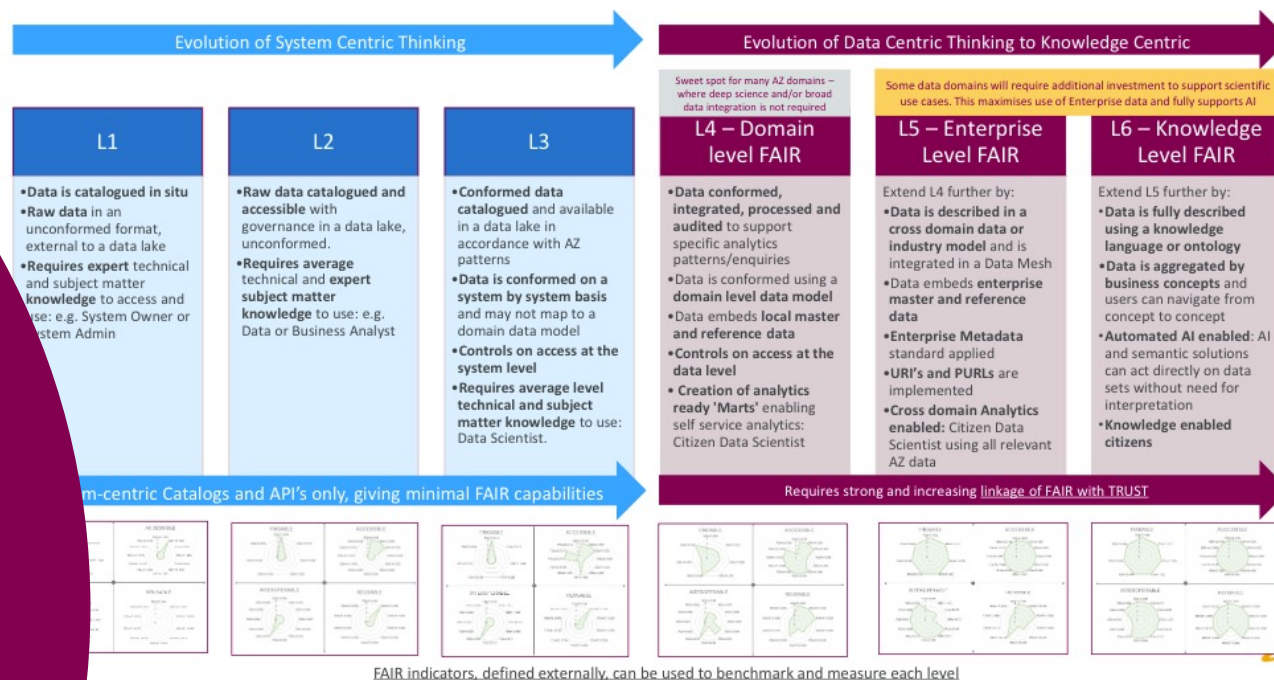


# What do we need to achieve FAIR?:

## The need for structured data and its consequences



# AZ Enterprise FAIR levels



ID	PRINCIPLE	INDICATOR_ID	INDICATORS	PRIORITY	METRIC	VIZ	S
1	F	P1	RDA-F1-01M Metadata is identified by a persistent identifier	Essential	1 – not being considered this yet	1	0
2		P1	RDA-F1-01D Data is identified by a persistent identifier	Essential	2 – under consideration or in planning phase	2	0
3		P1	RDA-F1-02M Metadata is identified by a globally unique identifier	Essential	1 – not being considered this yet	1	0
4		P1	RDA-F1-02D Data is identified by a globally unique identifier	Essential	2 – under consideration or in planning phase	2	0
5		P2	RDA-F2-01M Rich metadata is provided to allow discovery	Essential	1 – not being considered this yet	1	0
6		P3	RDA-F3-01M Metadata includes the identifier for the data	Essential	1 – not being considered this yet	1	0
7	A	P4	RDA-F4-01M Metadata is offered in such a way that it can be harvested and indexed	Essential	1 – not being considered this yet	1	0
8		A1	RDA-A1-01M Metadata contains information to enable the user to get access to the data	Important	1 – not being considered this yet	1	0
9		A1	RDA-A1-02M Metadata can be accessed manually (i.e. with human intervention)	Essential	3 – in implementation phase	3	0
10		A1	RDA-A1-02D Data can be accessed manually (i.e. with human intervention)	Essential	3 – in implementation phase	3	0
11		A1	RDA-A1-03M Metadata identifier resolves to a metadata record	Essential	1 – not being considered this yet	1	0
12		A1	RDA-A1-03D Data identifier resolves to a digital object	Essential	3 – in implementation phase	3	0
13	I	A1	RDA-A1-04M Metadata is accessed through standardised protocol	Essential	1 – not being considered this yet	1	0
14		A1	RDA-A1-04D Data is accessible through standardised protocol	Essential	3 – in implementation phase	3	0
15		A1	RDA-A1-05D Data can be accessed automatically (i.e. by a computer program)	Important	3 – in implementation phase	3	0
16		A1.1	RDA-A1-10M Metadata is accessible through a free access protocol	Essential	1 – not being considered this yet	1	0
17		A1.1	RDA-A1-10D Data is accessible through a free access protocol	Important	3 – in implementation phase	3	0
18		A1.2	RDA-A1-10D Data is accessible through an access protocol that supports authentication and authorisation	Useful	3 – in implementation phase	3	0
19	R	A2	RDA-A2-01M Metadata is guaranteed to remain available after data is no longer available	Essential	1 – not being considered this yet	1	0
20		I1	RDA-I1-01M Metadata uses knowledge representation expressed in standardised format	Important	1 – not being considered this yet	1	0
21		I1	RDA-I1-01D Data uses knowledge representation expressed in standardised format	Important	1 – not being considered this yet	1	0
22		I1	RDA-I1-02M Metadata uses machine-understandable knowledge representation	Important	1 – not being considered this yet	1	0
23		I1	RDA-I1-02D Data uses machine-understandable knowledge representation	Important	1 – not being considered this yet	1	0
24		I2	RDA-I2-01M Metadata uses FAIR-compliant vocabularies	Important	1 – not being considered this yet	1	0
25	I	I2	RDA-I2-01D Data uses FAIR-compliant vocabularies	Important	2 – under consideration or in planning phase	2	0
26		I3	RDA-I3-01M Metadata includes references to other metadata	Important	1 – not being considered this yet	1	0
27		I3	RDA-I3-01D Data includes references to other data	Useful	3 – in implementation phase	3	0
28		I3	RDA-I3-02M Metadata includes references to other data	Useful	1 – not being considered this yet	1	0
29		I3	RDA-I3-02D Data includes qualified references to other data	Useful	3 – in implementation phase	3	0
30		I3	RDA-I3-03M Metadata includes qualified references to other metadata	Important	1 – not being considered this yet	1	0
31	R	I3	RDA-I3-03D Metadata include qualified references to other data	Useful	1 – not being considered this yet	1	0
32		R1	RDA-R1-01M Plurality of accurate and relevant attributes are provided to allow reuse	Essential	3 – in implementation phase	3	0
33		R1.1	RDA-R1-10M Metadata includes information about the license under which the data can be reused	Essential	1 – not being considered this yet	1	0
34		R1.1	RDA-R1-10M Metadata refers to a standard reuse licence	Important	1 – not being considered this yet	1	0
35		R1.1	RDA-R1-10M Metadata refers to a machine-understandable reuse licence	Important	1 – not being considered this yet	1	0
36		R1.2	RDA-R1-20M Metadata includes provenance information according to community specific standards	Important	1 – not being considered this yet	1	0
37	R	R1.2	RDA-R1-20M Metadata includes provenance information according to a cross-community language	Useful	1 – not being considered this yet	1	0
38		R1.3	RDA-R1-30M Metadata complies with a community standard	Essential	1 – not being considered this yet	1	0
39		R1.3	RDA-R1-30D Data complies with a community standard	Essential	3 – in implementation phase	3	0
40		R1.3	RDA-R1-30M Metadata is expressed in compliance with a machine-understandable community standard	Essential	1 – not being considered this yet	1	0
41		R1.3	RDA-R1-30D Data is expressed in compliance with a machine-understandable community standard	Important	1 – not being considered this yet	1	0





# Enterprise FAIR standard: Defining AZ FAIR levels

Evolution of System Centric Thinking

Evolution of Data Centric Thinking to Knowledge Centric

L1

- **Data is catalogued in situ**
- Raw data in an unconformed format, external to a data lake
- Requires expert technical and subject matter knowledge to access and use

L2

- **Raw data catalogued** and accessible with governance in a data lake, unconformed.
- Requires average technical and expert subject matter knowledge to use

L3

- Conformed data catalogued and available in a data lake in accordance with AZ patterns
- Data is conformed on a system by system basis and **may not map to a domain data model**
- **Controls on access at the system level**
- Requires average level technical and subject matter knowledge to use: Data Scientist.

L4 – Domain level FAIR

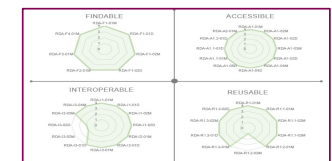
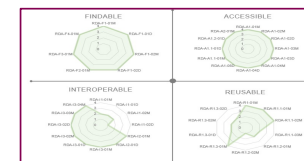
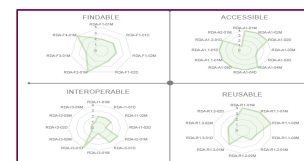
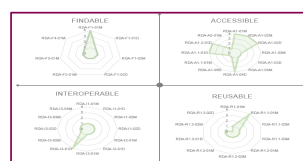
- Data conformed, integrated, processed and audited to support specific analytics patterns/enquiries
- **Data is conformed using a domain level data model**
- **Data embeds local master and reference data**
- Controls on access at the data level
- Creation of analytics ready 'Marts' enabling self service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

L5 – Enterprise Level FAIR

- Extend L4 further by:
- **Data is described in a cross domain data or industry model** and is integrated in a Data Mesh
  - **Data embeds enterprise master and reference data**
  - Enterprise Metadata standard applied
  - **URI's and PURLs are implemented**
  - Cross domain Analytics enabled: Citizen Data Scientist using all relevant AZ data

L6 – Knowledge Level FAIR

- Extend L5 further by:
- Data is fully described using a **knowledge language or ontology**
  - Data is aggregated by business concepts and users can navigate from concept to concept
  - Automated AI enabled: AI and semantic solutions can **act directly on data sets without need for interpretation**
  - **Knowledge enabled citizens**



Adapted from Colin Wood

FAIR indicators, defined externally, can be used to benchmark and measure each level



# How is the assessment carried out?

## Focus on Level 4

### Criteria for achieving Level 4 - Domain Level FAIR

Programmatic Data Provisioning	4. Is there programmatic provisioning, for example a data level API to access data? (ACCESSIBLE)	
	<ul style="list-style-type: none"><li>Is there programmatic provisioning to access the data?</li><li>Can data be accessed at a grain that allows consumption by grouping (e.g. by study (or equivalent grouping)) as well as bulk level basis?</li><li>Are outputs provided in a commonly used format?</li><li>Does the programmatic access mechanism provide metadata at entity and attribute level of the data it exposes?</li><li>If API is available, is it registered in Anypoint Exchange so access control can be coordinated via Mulesoft DeOps?</li></ul>	
Access Control	5. Are data product centric access controls in place to manage access to the data? (ACCESSIBLE)	
	<ul style="list-style-type: none"><li>Are the access control determined by metadata that defines access roles?</li><li>Is there a system or process in place that allows the data owner to revoke access if necessary?</li><li>Is there a data access policy from Data Office in place for the data / data set in question?</li><li>Is the access control orchestrated by a workflow tool?</li><li>Are the access control at an appropriate access grain and tied into metadata about critical data as defined by Note: Permission to access this data is independent of the system the data may have originated from.</li></ul>	
Reference and Master Data Identification	6. Has the data model (minimum metadata or conceptual/logical) undergone an IA and Data Office review (UI and reference data sources)? (INTEROPERABLE)	
	<ul style="list-style-type: none"><li>Has the data model for the data product undergone an IA review through IARB?</li><li>Has Master and Reference Data been identified?</li><li>Has Master and Reference Data sources been confirmed with Data Office (e.g. through IARB review)?</li></ul> <p>Note: Depending on the solution, conceptual or logical level might need to be reviewed. Level to determine "In future will need to have defined minimum/critical Master &amp; Reference data per domain"</p>	
RDM Service	7. Is the Platform consuming and/or publishing relevant reference data terms from a reference data solution	
	<ul style="list-style-type: none"><li>Are all domain related entities included?</li><li>Is the reference data in scope aligned to a standard vocabulary? / To what degree is the reference data in the vocabulary available in C/Ntne or will it be requested to be added to C/Ntne?</li></ul> <p>Note:</p> <ul style="list-style-type: none"><li>Interim tactical solutions do not count, unless it is approved as a standards by the Data Standards Councils</li><li>Accountability for Completeness is with the Data Office, as part of TRUST Framework. The Data Office prio may specify expected measures for completeness for current and legacy data.</li></ul>	

### Criteria for achieving Level 4 - Domain Level FAIR

NDM Service	8. Is the Platform consuming and publishing identifiers from domain specific Master Data Solutions? (INTEROPERABLE)	
	<ul style="list-style-type: none"><li>Is the platform/data product consuming identifiers from domain specific Master Data Solution?</li><li>Or is the platform/data product sourcing the master data in scope from the appropriate system of record?</li></ul> <p>Note: Interim tactical solutions do not count, unless it is approved as a standards by the Data Standards Councils.</p>	
Meta Data	9a. Is there a description of the minimum metadata terms available within Collibra? * (REUSABLE)	
	9b. OR Have business terms, defining tables and entities, been published to Collibra using the logical data model as a source? (REUSABLE)	
Data Audit	10. Are audit logs for the platforms published e.g. by sending logs to a suitable logging tool such as Splunk? (REUSABLE)	
	<p>This is not specifically aligned to FAIR practices, but is a requirement for GdP and should align to <a href="#">MHRA's GXP data integrity guide published</a> - <a href="#">MHRA Inspectorate (hmg.gov.uk)</a></p> <p>The audit trail is to provide secure recording of life-cycle details - creation, additions, deletion or alteration - of information in a record.</p> <p>The purpose of this is to embed TRUST in the data, but primarily to provide evidence of the journey of the data to allow it to be used with confidence. This should include:</p> <ul style="list-style-type: none"><li>Where the data originated</li><li>Loading information (Status and Time)</li><li>Data consumption: who or what utilised/accessed the data</li><li>Any alteration of the data</li></ul>	
Data Lineage	11. Stretch: Has upstream and downstream Data Lineage metadata been published to Collibra, identifying sources and consumers of data? (REUSABLE)	
	<ul style="list-style-type: none"><li>To what degree/coverage is data lineage available in Collibra for the dataset/data product in question?</li></ul>	

ID	PRINCIPLE	INDICATOR_ID	INDICATORS	PRIORITY	METRIC	VIZ	IS
1	F1	RDA-F1-01M	Metadata is identified by a persistent identifier	Essential	1 - not being considered this yet	1	0
2	F1	RDA-F1-01D	Data is identified by a persistent identifier	Essential	2 - under consideration or in planning phase	2	0
3	F1	RDA-F1-02M	Metadata is identified by a globally unique identifier	Essential	1 - not being considered this yet	1	0
4	F1	RDA-F1-02D	Data is identified by a globally unique identifier	Essential	2 - under consideration or in planning phase	2	0
5	F2	RDA-F2-01M	Rich metadata is provided to allow discovery	Essential	1 - not being considered this yet	1	0
6	F3	RDA-F3-01M	Metadata includes the identifier for the data	Essential	1 - not being considered this yet	1	0
7	F4	RDA-F4-01M	Metadata is offered in such a way that it can be harvested and indexed	Essential	1 - not being considered this yet	1	0
8	A1	RDA-A1-01M	Metadata contains information to enable the user to get access to the data	Important	1 - not being considered this yet	1	0
9	A1	RDA-A1-02M	Metadata can be accessed manually (i.e. with human intervention)	Essential	3 - in implementation phase	3	0
10	A1	RDA-A1-02D	Data can be accessed manually (i.e. with human intervention)	Essential	3 - in implementation phase	3	0
11	A1	RDA-A1-03M	Data identifier resolves to a metadata record	Essential	1 - not being considered this yet	1	0
12	A1	RDA-A1-03D	Data identifier resolves to a digital object	Essential	3 - in implementation phase	3	0
13	A1	RDA-A1-04M	Metadata is accessible through standardised protocol	Essential	1 - not being considered this yet	1	0
14	A1	RDA-A1-04D	Data is accessible through standardised protocol	Essential	3 - in implementation phase	3	0
15	A1	RDA-A1-05D	Data can be accessed automatically (i.e. by a computer program)	Important	3 - in implementation phase	3	0
16	A1.1	RDA-A1.1-01M	Metadata is accessible through a free access protocol	Essential	1 - not being considered this yet	1	0
17	A1.1	RDA-A1.1-01D	Data is accessible through a free access protocol	Important	3 - in implementation phase	3	0
18	A1.2	RDA-A1.2-01D	Data is accessible through an access protocol that supports authentication and authorisation	Useful	3 - in implementation phase	3	0
19	A2	RDA-A2-01M	Metadata is guaranteed to remain available after data is no longer available	Essential	1 - not being considered this yet	1	0
20	I1	RDA-I1-01M	Metadata uses knowledge representation expressed in standardised format	Important	1 - not being considered this yet	1	0
21	I1	RDA-I1-01D	Data uses knowledge representation expressed in standardised format	Important	1 - not being considered this yet	1	0
22	I1	RDA-I1-02M	Metadata uses machine-understandable knowledge representation	Important	1 - not being considered this yet	1	0
23	I1	RDA-I1-02D	Data uses machine-understandable knowledge representation	Important	1 - not being considered this yet	1	0
24	I2	RDA-I2-01M	Metadata uses FAIR-compliant vocabularies	Important	1 - not being considered this yet	1	0
25	I2	RDA-I2-01D	Data uses FAIR-compliant vocabularies	Useful	2 - under consideration or in planning phase	2	0
26	I3	RDA-I3-01M	Metadata includes references to other metadata	Useful	1 - not being considered this yet	1	0
27	I3	RDA-I3-01D	Data includes references to other data	Useful	3 - in implementation phase	3	0
28	I3	RDA-I3-02M	Metadata includes references to other data	Useful	1 - not being considered this yet	1	0
29	I3	RDA-I3-02D	Data includes qualified references to other data	Useful	3 - in implementation phase	3	0
30	I3	RDA-I3-03M	Metadata includes qualified references to other metadata	Important	1 - not being considered this yet	1	0
31	I3	RDA-I3-04M	Metadata includes qualified references to other data	Useful	1 - not being considered this yet	1	0
32	R1	RDA-R1-01M	Plurality of accurate and relevant attributes are provided to allow reuse	Essential	3 - in implementation phase	3	0
33	R1.1	RDA-R1.1-01M	Metadata includes information about the licence under which the data can be reused	Essential	1 - not being considered this yet	1	0
34	R1.1	RDA-R1.1-01M	Metadata refers to a standard reuse licence	Important	1 - not being considered this yet	1	0
35	R1.1	RDA-R1.1-01M	Metadata refers to a machine-understandable reuse licence	Important	1 - not being considered this yet	1	0
36	R1.2	RDA-R1.2-01M	Metadata includes provenance information according to community-specific standards	Important	1 - not being considered this yet	1	0
37	R1.2	RDA-R1.2-01M	Metadata includes provenance information according to a cross-community language	Useful	3 - in implementation phase	3	0
38	R1.3	RDA-R1.3-01M	Metadata complies with a community standard	Essential	1 - not being considered this yet	1	0
39	R1.3	RDA-R1.3-01D	Data complies with a community standard	Essential	3 - in implementation phase	3	0
40	R1.3	RDA-R1.3-02M	Metadata is expressed in compliance with a machine-understandable community standard	Essential	1 - not being considered this yet	1	0
41	R1.3	RDA-R1.3-02D	Data is expressed in compliance with a machine-understandable community standard	Important	1 - not being considered this yet	1	0

Solution	HBS Data Hub	
Service Now Link	TBD	
Blueprint Link	<a href="#">HBS Data Hub Solution Blueprint</a>	
ARB Review Date	28-Sep-22	
FAIR Assessor	Colin Wood	
FAIR Criteria	Score	Comments
1. Lake	1	Solution is fully integrated with Data Lake
2. Catalogued	0	Solution does not use data catalogues or metadata repositories
3. Critical Data	0	Solution and metadata not catalogued, so it is not possible to identify this.
4. Confirmation	1	Conceptual and Logical Data Models reviewed and approved at IARB
5. Programmatic Data Provisioning	1	API, delivered via Snowflake
6. Access Control	1	Full access control in scope for the blueprint
7. Reference and Master Data Identification	1	Full list of master and reference data in the blueprint and in this workbook
8. RDM Service	0.38	No use of RDM solutions - using data directly from operational sources
9. MDM Service	0.70	Using best available master data sources where available.
10. Metadata Vocabulary	0	Solution does not use data catalogues or metadata repositories
11. Data Audit	1	Auditing is built into the solution
12. Data Lineage	0	No use of metadata repositories to capture lineage
LEVEL 4 FAIR Score	64.37%	
Note that target score is 100% - but 120% is possible by embedding PURL Identifiers		

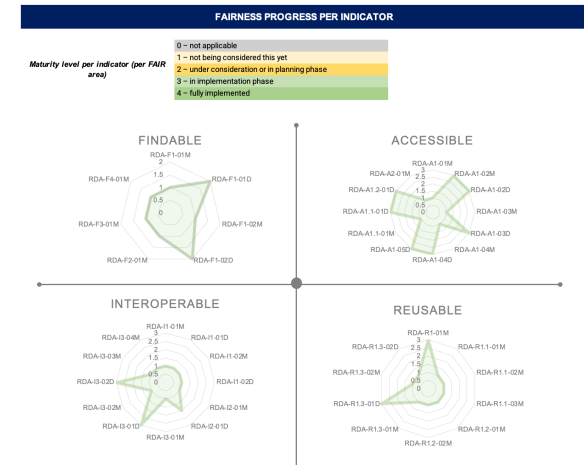
We recommend a separate access control review later in the project if

Build in support for ontologies sourced from CENTREE

More detailed plan for master data integration required

Catalogue in Collibra and publish the logical data model.

Publish lineage, specifying sources of data, into Collibra



# FAIRe(nough) Framework





# How FAIR is my data: the FAIRe(nough) benchmark at AstraZeneca

Ensure we target resources where the value is using Use case/business value as a driver

## Prioritised

### Use case example

As a ML/statistics scientist,  
I want to be able to build  
a **Drug Sensitivity Predictive Model**  
for **Disease X** patient  
treated with **drug Y**  
using **gene expression profiles**



Use case example (Data Science)

01

## Concepts and Data sources discovery

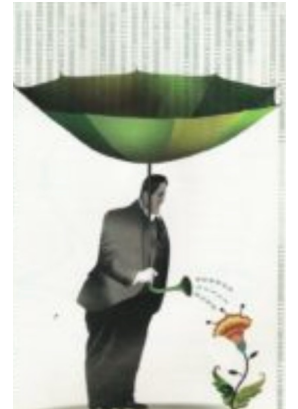
Find Patient level data sources where **indication X**, **response** following a **drug treatment X** and **OMIC data** are available

02

## Source and key elements gathering

Data Sources - High level metrics

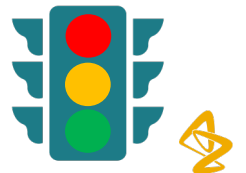
FAIRe Data – low level metrics



03

## FAIRe implementation

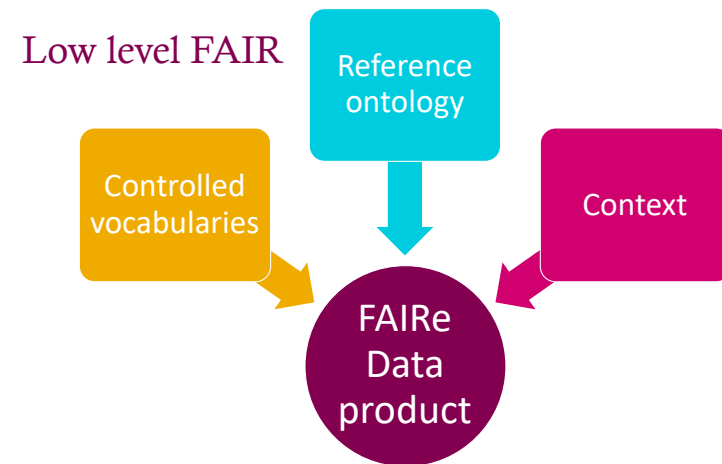
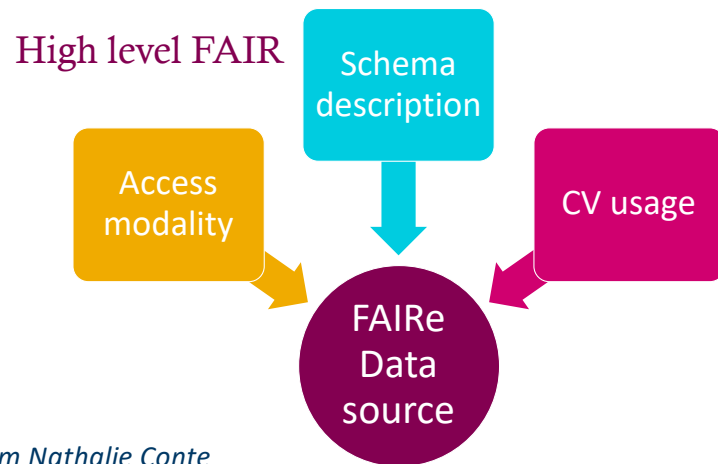
- FAIRe Data assessment (High/low level)
- Compute result and translate to impact
- FAIRe driven action:
  - Send back
  - Curate/ingest/buy



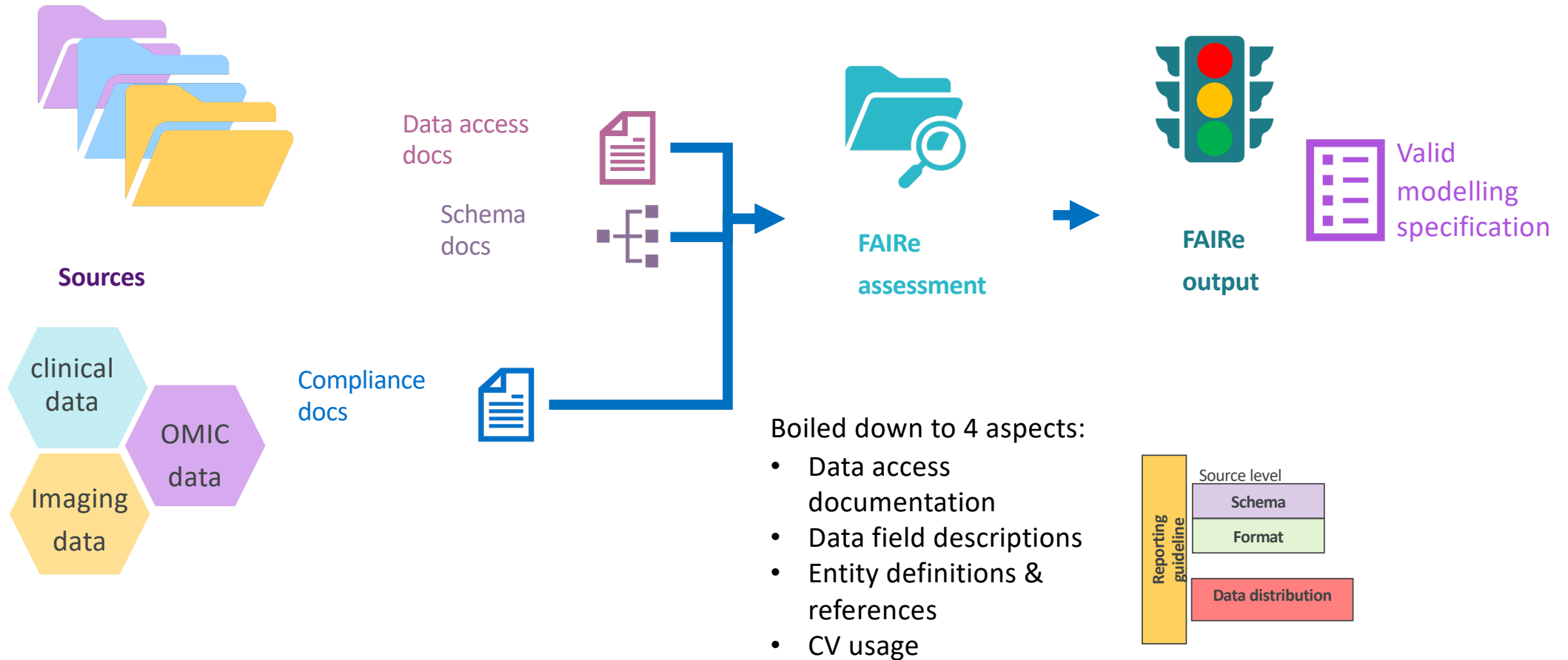
# FAIRe: Producing FAIRe(nough) data products

## Two dimensions of FAIRness ...

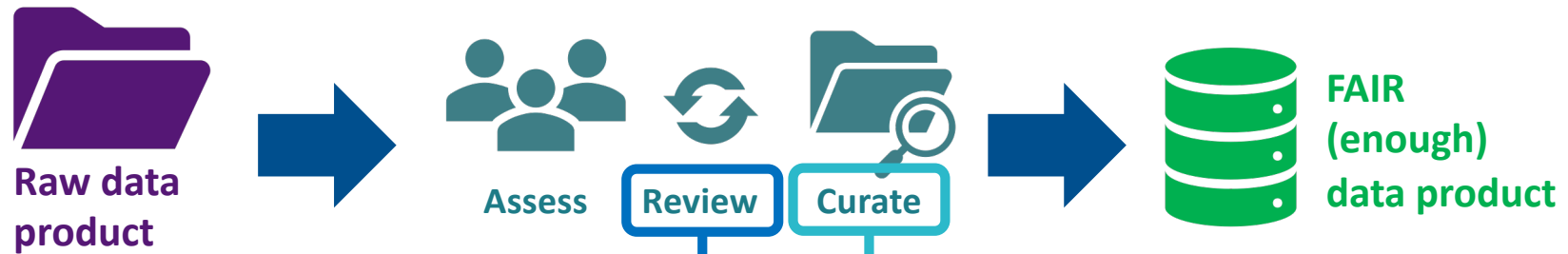
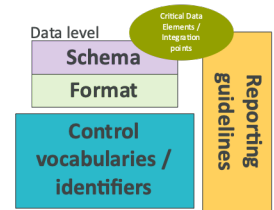
High-level FAIRness	Low-level FAIRness: FAIRe for a purpose
FAIRness of the data resource itself.	FAIRness of the data.
Does not require a data sample to calculate.	Requires a data sample to calculate.
Based on aspects of resource documentation, provisioning, access methods etc. (metadata attributes).	Based on metrics calculated from the data itself.
Fairly easy / cheap to define and calculate.	Less easy to define / calculate.
Calculated per resource.	Calculated via data assessment tools for a defined set of critical data elements in the context of a use case within/across data resources Completeness/adherence to Std/CV



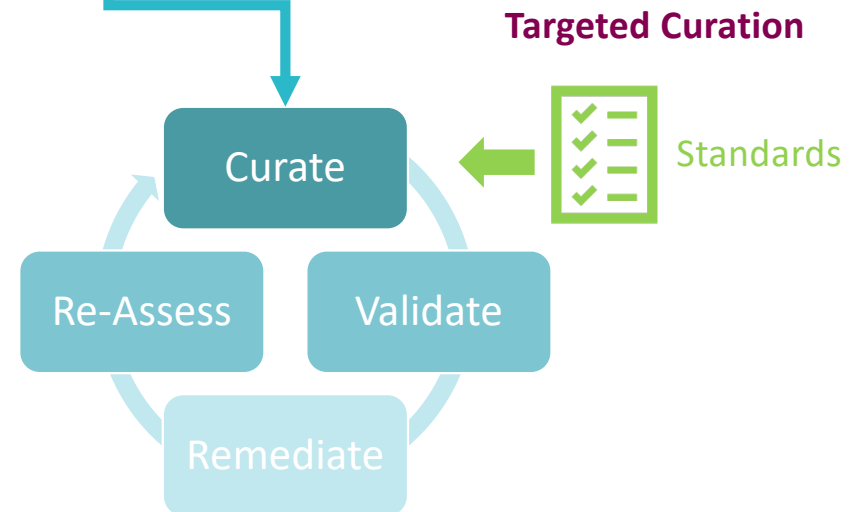
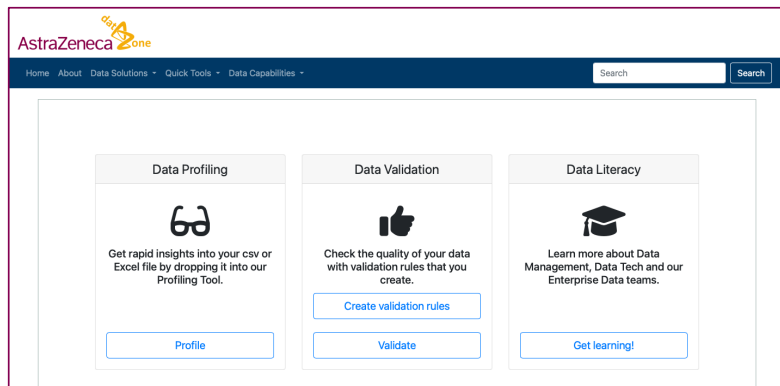
# High-level FAIR assessment



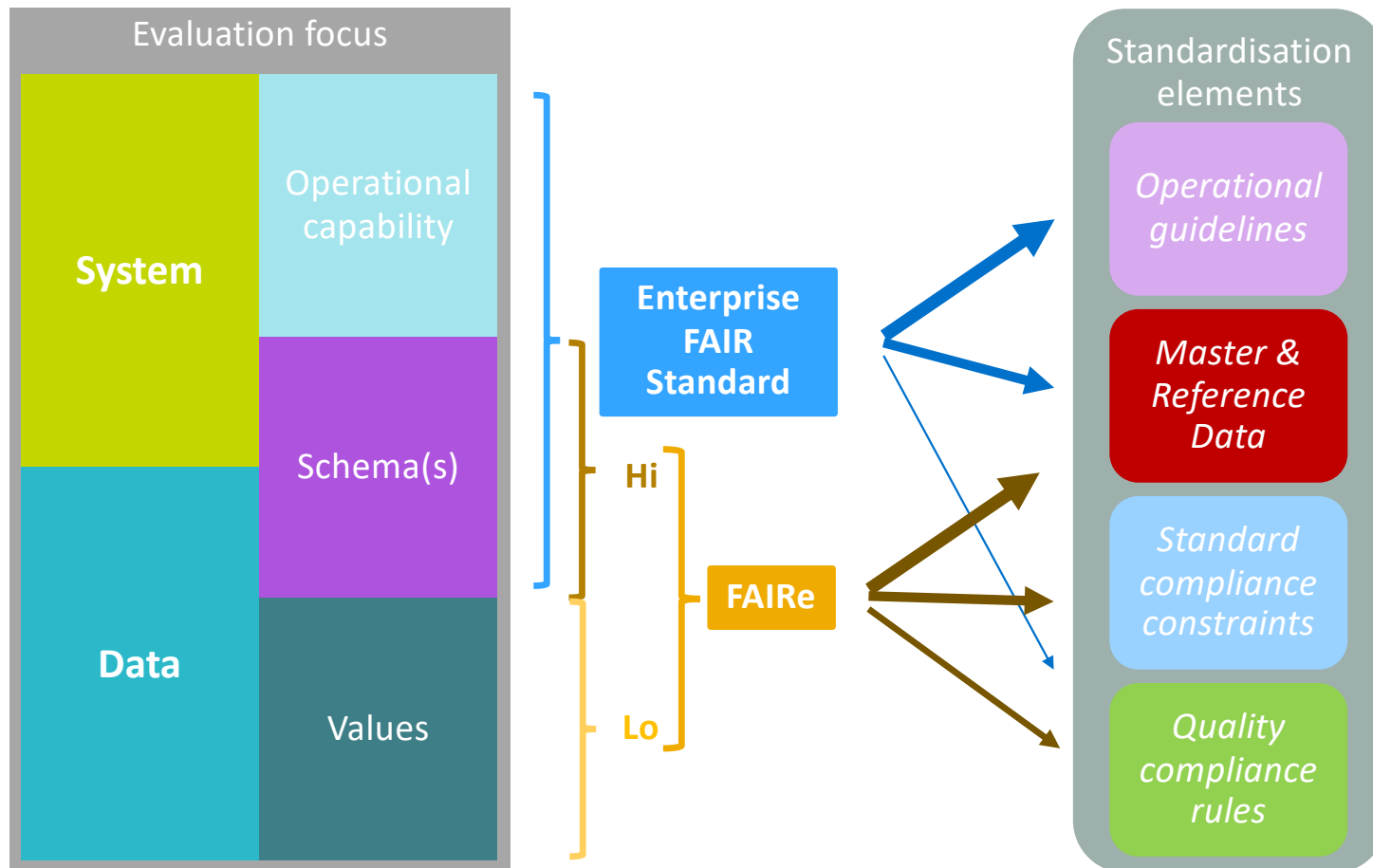
# Low-level FAIR assessment – Data assessment / curation loop



## AZ Data Zone tool



# Benchmarks scope





# Benchmarks assessment

## Enterprise FAIR Standard

- **Test cases:** Two independent system-level evaluations
- **Lessons learned:**
  - Defined as a **high-level** assessment benchmark
  - **Easy to measure**
  - **Quantitative output** for qualitative yes-no questions
  - **In scope:** system operational capability, data schema
  - **Out of scope:** data value level, qualitative/validation dimension of the data

## FAIRe(nough)

- **Test cases:** Data asset records in a single application
- **Lessons learned**
  - Current form is **integration specs guideline** rather than assessment benchmark
  - Implementation as a benchmark **challenging** at both high and low levels
  - **High level:**
    - Most elements can be met with Enterprise FAIR Standard
  - **Low level:**
    - Formalised most FAIRe questions as **validation rules** in DataZone Validator
    - Turning guidelines into quantitative scores is a challenge

## In practice



### High-level FAIR

- Work in progress to get a single score per level
- **North star** for system-level FAIR

### Low-level FAIR

- Generation of validation rule set
- Assessment of data against **validation rules**
- Governance of individual **data elements** is key



# Next steps



Global Strategic partnership

- Mature development of **low-level assessment** into **defined set of rules**
  - Consolidation of **Core Data Elements** reference list
  - Representation and implementation of **business rules** around Core Data Elements into validation toolkit
  - Mature into **common validation toolkit**
- Mature **high-level assessment**
  - **Consolidation of alignment** between FAIR high-level information standard and Enterprise FAIR levels
- FAIR **measures to be computable, cohesive and transparent** to all Data users
  - FAIR **toolkit** development
  - **Governance & policies** alignment
  - Focus on **sustainability**
- **Improve/publish the framework** through iteration & strategic alignment
- **Open for Collaboration** with relevant initiatives! Please contact us 😊



Thanks for  
listening!



**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](http://www.astrazeneca.com)

