



June, 2026

Data Director Agentic AI Tool Blueprint

Open, Technology-Agnostic Blueprint v0.1

Formal Recommendation of the
RDA Data Director Agentic AI Blueprint Working Group

Part of the RDA-Microsoft Global Agentic AI Initiative

Programme Lead: Connie Clare (RDA)

Programme Facilitators: Connie Clare (RDA), George Earl (Microsoft),
Trish Radotic (RDA) and Benjamin Wright-Jones (Microsoft)



Data Director

Agentic AI Tool Blueprint

Version 0.1 - DRAFT

Initiative: RDA-Microsoft Global Agentic AI Initiative

Status: First Draft for RDA Community Review

Version: 0.1

Date: June 2026

Programme Lead: Connie Clare (RDA)

Programme Facilitators: Connie Clare (RDA), George Earl (Microsoft), Trish Radotic (RDA), Benjamin Wright-Jones (Microsoft)

Blueprint Authors: See [Appendix C: Acknowledgements](#)

This document is a community-developed technical specification for the Data Director as an open, technology-agnostic agentic AI tool. It is not a finished product or piece of software. It is freely accessible, transparently developed and governed in line with RDA Guiding Principles.¹ It defines what the Data Director should do and how it should behave, without prescribing any specific AI model, cloud provider or infrastructure stack.

¹ <https://www.rd-alliance.org/about/>

Table of Contents

1. Executive Summary	7
2. Introduction	8
2.1 About the Working Group	8
2.2 Purpose and Scope	8
2.3 What is an Open, Technology-Agnostic Blueprint?	9
2.4 The Data Director: Tool Definition	9
2.5 Intended Users	10
3. Problem Definition	10
3.1 The Research Data Challenge	10
3.1.1 Core Gaps: Directly addressed by the Data Director	11
3.1.2 Boundary Conditions: Gaps that shape how the Data Director should behave	13
3.1.3 Acknowledged Gaps: Not directly addressed by the Data Director	14
3.2 Vision and Goals	15
3.3 Out of Scope	16
4. Research Context	16
4.1 Uses, Users and Expected Outcomes	16
5. Blueprint Considerations	26
5.1 Governance and Roles	26
5.2 Legal, Ethical and Sovereignty Prerequisites	27
5.3 Resolving Conflicting Requirements	27
5.4 AI Behaviour and Automation	27
5.5 Operational Readiness	28
5.6 Costs, Sustainability and Evaluation	29
6. Architecture Principles	29
7. Functional Specification	34
7.1 Functional Requirements	34
7.2 Non-functional Requirements	43
8. Process Flow	49
8.1 How to read the process flow	50
8.2 Six-Phase Process Flow for Research Data Publication	50
9. Reference Architecture	58



9.1 Architecture Overview	58
9.2 Presentation Layer	59
9.3 Application Layer	60
9.4 AI Layer	60
9.5 Data Layer	61
9.6 Platform and Integration Layer	62
10. Evaluation Criteria	63
10.1 Blueprint Quality Criteria	64
10.2 Evaluation Process	69
10.3 Future Evaluation Considerations	69
11. Implementation Guidance	70
11.1 Cross-cutting Practicalities	71
11.2 Research Institutions	71
11.3 National Research Infrastructures	72
11.4 Commercial or Non-profit Vendors	73
12. Blueprint Governance	73
12.1 Ownership and Stewardship	73
12.2 Maintenance and Continuity	73
12.3 Platform and Transparency	74
12.4 Review and Evolution	74
12.5 Scope and Sustainability	74
13. Open Questions and Future Iterations	75
13.1 Unresolved Design Questions	75
13.2 Roadmap for Future Iterations	76
14. Appendices	78
Appendix A: Glossary of Terms	78
Appendix B: Related Resources	81
Appendix C: Acknowledgements	89
Blueprint Authors	89
Appendix D: Potential Data Director Blueprint Implementors	93
Appendix E: Tool Usage	94
Appendix F: Partners and Document Information	94



About the RDA	94
About Microsoft	95
Licence and Contact	95



Document Control

Version History

Version	Date	Change Summary
0.1	June 2026	Initial draft for community review

1. Executive Summary

This Blueprint sets out a community-developed, technology-agnostic specification for the Data Director, a proposed agentic AI tool to support researchers in preparing, documenting and depositing research data for publication in alignment with the FAIR principles. It is the formal output of the Data Director Agentic AI Working Group, established through an RDA-Microsoft global community consultation on agentic AI in research, and represents version 0.1: a first draft for community review rather than a finished product.

The case for the Data Director rests on six interconnected gaps in current research data practice. Four core gaps, relating to limitations in researcher capacity, inconsistent metadata standards, uneven access to data support, and uncertainty navigating funder, journal and institutional policies, are directly addressed by the tool's functional specification. Two further gaps, concerning the quality and provenance of data and metadata, and the governance of sensitive, ethically or culturally significant data, including Indigenous data under the CARE principles, establish boundary conditions for the tool's behaviour. Two acknowledged gaps, around infrastructure sustainability and the recognition of data outputs in academic reward systems, fall outside the Data Director's scope but are documented for completeness.

This problem definition is translated into 14 architecture principles spanning FAIR and CARE alignment, openness, security, human oversight, traceability, interoperability, accessibility, sovereignty and environmental impact. These are operationalised through 12 functional requirements, covering repository recommendation, FAIR-aligned metadata generation, persistent identifier support, Data Management Plan verification, metadata remediation, provenance tracking, and discovery and linking, alongside 17 non-functional requirements defining the quality standards any implementation must meet. A six-phase process flow and vendor-agnostic reference architecture show how these come together in practice, from project planning through to monitoring and reuse of research data after publication.

With no reference implementation yet existing, the Blueprint sets out forward-looking evaluation criteria, implementation guidance for research institutions, national research infrastructures and commercial or non-profit vendors, and a governance model proposing the RDA as steward of the specification.

Finally, the Blueprint is transparent about what remains unresolved. Its open design questions and roadmap for future Blueprint iterations should be read not as gaps in the work, but rather as a sign that the eight-week consultation has surfaced genuine tensions and future considerations. The breadth of these questions reflects the depth of community engagement and offers a basis for ongoing discussion as the Blueprint is developed further.

2. Introduction

2.1 About the Working Group

The Data Director Blueprint has been produced by the [Data Director Agentic AI Working Group](#),² established following [a global community consultation on agentic AI in research](#)³ conducted by the [Research Data Alliance \(RDA\)](#)⁴ in collaboration with [Microsoft](#)⁵ in November 2025. That consultation built on prior engagement Microsoft undertook with around 50 research institutions across the UK, which had identified data preparation and FAIR-aligned data sharing as a significant and widely felt pain point across the research lifecycle. The consultation explored current use of agentic AI by researchers and sought perspectives on its value across the research lifecycle, with agentic AI defined as '*artificial intelligence systems capable of autonomous operation with human oversight*'.

Three tools emerged as clear community priorities: the Literature Librarian, the Data Director, and the Funding Finder. Of these, the **Data Director**, designed to support research data preparation and sharing in alignment with the [FAIR principles](#),⁶ was identified as the tool most closely aligned with the RDA's mission, vision, and [guiding principles](#),⁷ placing data sharing, interoperability, and FAIR-aligned data management practices at its core. It was therefore selected as the focus for this Blueprint.

The Working Group operated over eight weeks through a community-driven process comprising two facilitated sessions per week, covering global time zones and open to all. The facilitation was supported by Microsoft within the framework of the [RDA public-private partnerships](#).⁸

2.2 Purpose and Scope

This Blueprint is a written functional and non-functional specification that defines what a Data Director agentic AI tool should do and how it should work. The development of any deployable implementation or prototype solution is out of scope for this phase of work and may be explored by the community in future. The Blueprint constitutes the Working Group's formal output under the [RDA Recommendations and Outputs process](#).⁹

² <https://www.rd-alliance.org/groups/data-director-agentic-ai-Blueprint/activity/>

³ https://www.rd-alliance.org/value_rda/rda-and-artificial-intelligence/global-community-priorities-for-agentic-ai-development/

⁴ <https://www.rd-alliance.org/>

⁵ <https://www.microsoft.com/>

⁶ <https://doi.org/10.1038/sdata.2016.18>

⁷ <https://www.rd-alliance.org/about/how-rda-works/>

⁸ <https://www.rd-alliance.org/wp-content/uploads/2024/03/RDA-Private-Sector-Engagement-Framework-v1.0.pdf>

⁹ <https://www.rd-alliance.org/recommendations-and-outputs/>

The Blueprint covers the following areas:

- The background and context for research data publication in alignment with the FAIR principles
- The problem definition and case for a Data Director agentic AI tool
- The architectural and design principles that should govern the tool
- A full functional specification, including a process flow and a high-level reference architecture diagram
- Evaluation criteria and implementation guidance
- Blueprint governance and open questions

The Blueprint is intended to be sufficiently detailed for a range of diverse implementers: research institutions, national research infrastructures, and commercial or non-profit vendors, to use as the basis for implementation. Existing tools and platforms may overlap with some capabilities described here. The Blueprint does not seek to replace or endorse any specific solution, but to provide a common, community-agreed specification against which tools can be developed or evaluated.

2.3 What is an Open, Technology-Agnostic Blueprint?

This Blueprint is a community-developed, openly governed specification, not a product or piece of software. It is designed to be implemented by any organisation regardless of their technology stack, and is defined by three principles:

- **Open:** Freely accessible, transparently developed, and community-governed in line with RDA's Guiding Principles and general principles of open science.
- **Technology-agnostic:** Defines what the tool should do and how it should behave, without prescribing which AI model, cloud provider, infrastructure, or data repository to use.
- **Living and versioned:** A specification designed to evolve through community review as needs, AI capabilities, and the broader research infrastructure landscape evolve. Version 0.1 represents the first community-endorsed draft.

2.4 The Data Director: Tool Definition

The Data Director is a proposed, technology-agnostic, human-supervised agentic AI tool that supports research data documentation, publication, sharing, and reuse. It does so by helping researchers create ethically and scientifically responsible metadata and data documentation, identify appropriate repositories, metadata standards and persistent identifier options, and meet relevant funder, institutional, disciplinary, and governance requirements.

As described in Section 2.1, the Data Director was identified as a community priority through the RDA-Microsoft global consultation and selected as the focus for this Blueprint based on its close alignment with the RDA's mission and guiding principles.

Community engagement identified three core capabilities as the starting point for the tool, listed here in the order in which they arise across the data publication workflow:

1. Recommending appropriate data repositories in line with funder, journal and institutional requirements
2. Automatically suggesting metadata in alignment with the FAIR principles
3. Supporting the assignment of persistent identifiers (PIDs) for datasets

These capabilities represent the community-validated foundation from which the Working Group's programme of work began. The Working Group developed them into the detailed functional specification set out in this Blueprint, expanding beyond the initial three capabilities to encompass a broader range of uses, functional requirements, and goals identified through the programme.

2.5 Intended Users

The Data Director is designed primarily for researchers, from early-career researchers to principal investigators, who need to prepare and share their data in alignment with the FAIR principles across the research lifecycle, but may lack the time, training, or institutional support to do so.

Data support professionals are a secondary but important user group, both as direct users of the tool and as those who may configure, deploy, or provide support for it on behalf of their institution.

While institutions, journals, and funding organisations are not direct users of the tool, they are key stakeholders whose requirements, particularly around data repository selection and policy requirements, shape what the tool recommends.

3. Problem Definition

3.1 The Research Data Challenge

The research community faces significant barriers to making data outputs FAIR (Findable, Accessible, Interoperable and Reusable). The following challenges were identified through the global community consultation on agentic AI in research conducted by the RDA and Microsoft in November 2025 and further developed and prioritised during the Blueprint programme.

Each is framed as a **gap**, a defined deficit between current practice and what is needed to share data effectively, providing a structured basis for the solutions and recommendations that follow. These gaps are wide-ranging and deeply embedded in research culture, infrastructure, and incentive structures. The Blueprint does not claim to address them in full; rather, it responds to the most common and pressing challenges identified by the community, providing a foundation that implementers can build on. The gaps are presented here to set the scene and establish the context for the Data Director's design, not to imply that a single tool can resolve them entirely. The gaps are presented here in order of community priority.

3.1.1 Core Gaps: Directly addressed by the Data Director

Capacity Gap: Researchers lack the time and expertise to manage data responsibly across the research lifecycle

A central challenge is that researchers lack both the training and the incentive to produce high-quality metadata. Throughout this Blueprint, metadata is understood in its full scope, encompassing not only discovery fields but variable definitions, data structure, provenance, quality assessments, access conditions, and applicable governance constraints. Metadata creation is not typically part of academic training or curricula, meaning many researchers do not know what good metadata looks like, how to produce it, or why it matters. They are not rewarded for investing time in it, and where a dataset is produced by a cross-disciplinary or cross-institutional team, priorities around metadata quality are not always shared. The result is that researchers tend to fulfil only the minimum requirements, with the burden of support falling on libraries and institutions. This reflects a deeper cultural issue: research data management is not yet treated as a core research responsibility, but as an administrative obligation to be minimised.

This challenge spans the full research lifecycle, not only the point of data publication. Decisions made early in a project, such as file formats, storage, and documentation practice, have significant consequences for whether data can ultimately be shared and reused, yet researchers are often unaware of these dependencies until problems arise. This is becoming more acute as AI systems are increasingly used to process research data, making metadata precision a prerequisite not only for discoverability but for scientific replicability and accountability.

A related challenge is that many researchers are not familiar with community-recognised repositories, standards registries, or vocabulary services, and lack a reliable basis for evaluating the trustworthiness of the resources they do encounter. Data Management Plans (DMPs) illustrate the broader problem well: researchers struggle to navigate competing institutional, funder, and collaborator requirements, to balance open science with research security and data protection, and to keep DMPs current as projects evolve.

Standards Gap: Inconsistency in metadata standards creates significant barriers to data sharing and reuse

There is no universally agreed interpretation of the FAIR principles, though the principles themselves provide a broadly accepted direction of travel that this Blueprint adopts as its orienting framework. What constitutes a complete metadata record varies significantly between communities and data repository types. This variability is compounded by the growth of interdisciplinary research, which requires researchers to navigate multiple, sometimes conflicting, standards frameworks simultaneously. The difficulty of gaining a reliable overview of the standards relevant to any given discipline compounds this further.

Where standard vocabularies or ontologies do not exist, these challenges deepen further. The harmonisation of heterogeneous metadata schemas, the conversion between them, and the consistent definition of variables, frequently overlooked but fundamental to interoperability,

remain largely unresolved and under-supported. Efforts to achieve standardisation have historically struggled, as community-specific needs often sit in tension with universal solutions.

These inconsistencies affect data discovery, making it difficult to locate relevant data outside one's own discipline, and carry equity implications: the costs of meeting evolving standards may disadvantage smaller repositories, risking the loss of important community resources.

Support Gap: Researchers lack access to skilled human support

Access to professional data support, such as data stewards or librarians, is unequally distributed across the research community. While smaller institutions are most visibly affected, the challenge is broader: large, devolved organisations face similar problems due to asymmetrical provision across faculties and disciplines. Even where support staff exist, they typically offer general guidance rather than hands-on assistance with tasks such as metadata creation. Access to timely support is a further constraint. Researchers frequently need help at short notice, and availability cannot be assumed.

Many researchers are unaware that support services exist at all, even at well-resourced institutions. Where support is available, it tends to concentrate around the point of data publication, with limited sustained curation effort afterwards, particularly as new standards emerge that may affect already-published datasets. A significant volume of research data has already been published without FAIR-aligned metadata, and the remediation of this existing material represents an important but largely unmet need.

Policy Gap: Researchers face uncertainty navigating divergent and sometimes conflicting data policies

The challenge here is broader than data repository selection alone. Researchers often lack clarity about what they are expected or required to do when publishing research data, and must navigate institutional, funder, and journal requirements that do not always align. Inconsistency in policy communication is as much an issue as policy conflict itself. Where institutions provide clear direction, the problem is significantly reduced.

Data repository selection is one practical manifestation of this challenge. The repository landscape is fragmented, leaving researchers without reliable criteria to guide their choices and forcing many to rely on general search engines as a starting point. Despite longstanding community efforts, consensus on what constitutes a good data repository has not been reached. The cost of data storage adds further complexity, as funding constraints can limit available options regardless of policy requirements. Journals also play a role here: many journals mandate that research data underpinning a journal article publication must be deposited in a specific repository or meet particular standards. These journal data policies must be considered alongside funder and institutional requirements when selecting a repository for data publication and can introduce additional or conflicting constraints.

The gap is best understood as one of comprehension: researchers struggle to understand the full range of requirements they are expected to meet, rather than facing simply a question of which data repository to use.

3.1.2 Boundary Conditions: Gaps that shape how the Data Director should behave

Quality and Provenance Gap: Researchers and tools must be able to demonstrate the trustworthiness and traceability of data and metadata

A core challenge is provenance: the need to know where a result comes from, how it was produced, and whether it can be trusted as the basis for new work. This concern is growing in urgency as increasing volumes of data and metadata are AI-generated, making traceability more critical and more difficult simultaneously. Any output from an agentic AI tool must include basic provenance information: the sources used, the model involved, the processing date, and a persistent identifier, as well as a record of whether human curation has taken place.

Closely related is the challenge of quality: both the quality of metadata, assessed against community standards and canonical semantics (a standard, universal representation of meaning and structure across datasets), and the quality of the data itself, which should be reflected as an explicit field within the metadata record. These two challenges are deeply interconnected: poor provenance undermines quality assessment, and poor quality metadata makes provenance harder to establish.

Provenance is already a core component of the FAIR principles, meaning tools that support FAIR-aligned metadata creation should produce appropriate provenance records as a matter of course. As AI plays a greater role in research data workflows, robust and standardised approaches to recording provenance and assessing quality become a fundamental requirement for research integrity.

Governance Gap: Researcher uncertainty about how to prepare data for sharing under ethical, legal, cultural or proprietary constraints

A significant challenge is researcher uncertainty about how to prepare data for sharing where ethical, legal, cultural, or commercial constraints apply. This uncertainty extends to data privacy, affecting not only researchers but data support professionals themselves. At its core, the gap reflects a lack of clarity about what is required and permissible when handling sensitive data, particularly around consent, confidentiality, and reuse. Cross-national research compounds this further, requiring researchers to navigate differing and sometimes conflicting requirements around data protection, rights holder obligations, and permitted reuse across jurisdictions.

The [CARE principles for Indigenous data](#)¹⁰ represent a specific and important dimension of this gap. Questions have been raised about whether agentic AI is compatible with CARE principles at all, and how alignment between the two might be achieved. Meaningful engagement with the CARE principles requires direct engagement with relevant Indigenous and community stakeholders and may involve interaction with external systems such as [Local](#)

¹⁰ <https://doi.org/10.5334/dsj-2020-043>

Contexts.¹¹ This remains a complex but important area for consideration in the development of any tool operating in this space.

Research data governance considerations are best addressed at the data management planning stage. However, the degree to which an AI tool can support governance decisions in a meaningful and trustworthy way is limited, particularly for edge cases that are difficult even for experienced human practitioners, and human oversight remains essential.

3.1.3 Acknowledged Gaps: Not directly addressed by the Data Director

Sustainability Gap: Difficulty funding and sustaining research data infrastructure at institutional scale

While the Data Director supports PID assignment as part of its core functionality, the broader challenge of funding and sustaining PID infrastructure at institutional scale falls outside its scope. In practice, the cost burden varies considerably; some repositories integrate PID assignment directly at no additional cost, and the significance of this as a barrier differs across institutions and contexts.

More pressing are the wider infrastructure challenges faced particularly by smaller institutions, where limited technical knowledge and restricted access to IT specialists make it difficult to automate even basic research data management processes. The longer-term sustainability of repositories hosting research data, including appropriate recognition and credit for those services, is also a concern that falls beyond what the Data Director can address.

There is a risk that intensive delegation to the agentic AI tools, such as the Data Director, could create dependency and hinder the development of in-house data management capacity over time. This is an important consideration for institutions deploying the tool.

Recognition Gap: Academic reward systems do not adequately recognise data outputs or data management efforts

Academic reward systems are weighted towards journal article publications, leaving research data management and sharing largely unrecognised. The immaturity of data citation metrics makes it difficult to demonstrate the impact of data sharing, creating a structural disincentive to investing in good research data management practice. While this gap cannot be addressed through tooling alone, reducing the effort required to share data well and providing researchers with metrics on how their data is being reused can lower barriers and offer a degree of recognition within the current system.

¹¹ <https://localcontexts.org/>

The **Sustainability** and **Recognition** Gaps are documented here to provide a complete picture of the research data landscape. They fall outside the scope of the Data Director in its current form but may be considered in future iterations or addressed by other tools and systems in development.

The four core gaps: **Capacity**, **Standards**, **Support**, and **Policy**, define the primary research data challenges the Data Director is designed to address. The **Quality and Provenance** and **Governance Gaps** establish the boundary conditions that shape how the tool should behave. Together, these six gaps provide the evidence base for the functional specification that follows.

3.2 Vision and Goals

The Data Director will support any researcher, regardless of institution, discipline, or level of data expertise, and the data support professionals who assist them, to prepare and deposit research data for publication in alignment with the FAIR principles, narrowing the gap between those with access to skilled data support and those without.

The following goals represent the primary starting point for the Data Director in this iteration. They are not exhaustive; further goals and capabilities are identified throughout the Blueprint and in the roadmap for future iterations in Section 13.2.

1. Simplify data repository selection

The Data Director will recommend appropriate data repositories based on funder, journal, and institutional requirements, so that researchers can make well-informed deposit decisions with confidence, regardless of their familiarity with the policy landscape.

2. Reduce the burden of research data documentation and metadata creation

The Data Director will enable researchers to select the most appropriate community and domain-relevant standards and schemas for their data, and automatically suggest metadata and supporting documentation, including data dictionaries and README files, so that researchers can produce FAIR-aligned records without relying on specialist data support.

3. Support the assignment of persistent identifiers (PIDs) for datasets

The Data Director will support the assignment of persistent identifiers (PIDs) for research datasets, so that data outputs are citable, discoverable, and persistently accessible regardless of the researcher's institutional resources.

4. Connect datasets to the broader research landscape

The Data Director will link datasets to related research outputs, including journal article publications, grants, software, and people, so that research data is discoverable in context and its contribution to the wider research record is visible.

3.3 Out of Scope

While the Data Director addresses several key challenges researchers face when depositing research data for publication, this Blueprint iteration has defined boundaries.

The Data Director does not:

- Support comprehensive data management across the full research data lifecycle; it is designed to support research data publication only.
- Replace human expertise, substitute for specialist research data management professionals, or assess the scientific validity or methodological soundness of research data.
- Take irreversible actions, including depositing research data for publication, without explicit human approval.
- Advise on how to meet legal, ethical, or governance requirements; it may flag relevant requirements, but researchers and institutions remain responsible for decisions that depend on local context and specialist knowledge. Note that the governance requirements shaping how the Data Director itself is designed and deployed are a separate matter, addressed in Section 5.
- Store, preserve, or provide repository infrastructure; it supports the preparation and submission of data for deposit but does not retain data itself.

The full boundary between what is included in this iteration and what may be reserved for future development is set out in the uses table in Section 4.1.

4. Research Context

4.1 Uses, Users and Expected Outcomes

This section sets out the range of uses the Data Director is designed to support, the users it serves, and the expected outcomes of each use. The primary outcome is to deliver the vision set out in Section 3.2: enabling any researcher to prepare and deposit research data for publication in alignment with the FAIR principles, regardless of their expertise or institutional support.

A **use** describes a specific situation in which a real user interacts with the tool to solve a real problem. Each use answers: who is doing what, why, and what happens as a result. The Working Group identified a range of uses for how the Data Director could support researchers and data professionals. Each has been assigned one of two statuses (Table 1):

Confirmed (green): in scope for this iteration of the Data Director. Confirmed uses are directly supported by a functional requirement in Section 7.1.

Future Consideration (yellow): addressing important gaps and challenges identified by the community, but beyond the scope of the current programme. These uses represent a 'community wish list' for future iterations of the Data Director. They are documented here to inform ongoing development; to signal the direction the community would like to see the tool

evolve, and to provide a foundation for future Blueprint versions or related agentic AI tools in development.

The volume of future considerations reflects the richness of community input rather than a lack of ambition. Given the eight-week timeframe of the Blueprint programme, not all identified uses could be fully discussed, specified, or progressed to a point where they could be confirmed for this iteration. All future considerations are nonetheless valued contributions and are strongly encouraged for follow-up in future iterations of the Data Director or through related community-led initiatives.

Table 1. Uses, users and expected outcomes of the Data Director, showing confirmed uses for this iteration and community-identified future considerations.

Use <i>A specific interaction between user and tool</i>	User(s) <i>Who initiates or benefits from this use</i>	Expected Outcome <i>What the tool produces or enables</i>	Aligned Gap(s) <i>The research challenge addressed in Section 3.1</i>	Status <i>In scope now, or identified for future development</i>
1. REPOSITORY SELECTION				
Search and recommend appropriate data repositories based on funder, institutional, journal, and disciplinary requirements	Researcher, Data Support Professional	A suitable repository is identified; conflicts between funder, institutional, journal and disciplinary requirements are flagged for human review; metadata recommendations for the selected repository are provided	Policy, Standards, Support	Confirmed
Set access conditions and assign a licence for data reuse, including open, restricted, and embargoed options	Researcher, Data Support Professional	Access conditions and licence are defined and attached to the dataset; open, restricted, and embargoed options are supported; funder or journal mandates are flagged where relevant	Policy, Governance	Future Consideration
Generate a private review link for manuscript peer reviewers to access an embargoed dataset before the associated journal article is published	Researcher, Data Support Professional	Peer reviewers can access the dataset under embargo without it being made publicly available; access is time-limited and auditable	Policy, Governance	Future Consideration

Monitor embargo status and alert the researcher when a dataset should be released to public access following journal article publication	Researcher, Data Support Professional	The researcher is notified when the embargo period has elapsed and the associated journal article has been published; release to public access is confirmed by the researcher before being enacted	Policy, Governance	Future Consideration
Check repository certification and trustworthiness, for example against CoreTrustSeal ¹² or FAIRsharing ¹³ criteria	Researcher, Data Support Professional	Repository trustworthiness is surfaced using recognised schemes such as CoreTrustSeal and FAIRsharing; free certification options are prioritised over paid alternatives	Policy, Support	Future Consideration
Provide guidance on managing and sharing very large data files, including file compression, transfer methods, and upload requirements for discipline-specific repositories	Researcher, Data Support Professional	Researchers are guided towards appropriate solutions for large file storage, transfer, and format conversion; barriers to sharing large datasets are reduced	Capacity, Support	Future Consideration
2. METADATA AND DOCUMENTATION				
Create FAIR-aligned and discipline-appropriate metadata, including controlled vocabularies, ontologies, and community standards such as CARE principles where applicable	Researcher, Data Support Professional, Institution	A draft metadata record is generated using appropriate controlled vocabularies (standardised term lists) and ontologies (structured concept frameworks); the researcher reviews and approves before submission	Capacity, Standards, Support	Confirmed
Suggest appropriate vocabularies and ontologies	Researcher, Data Support Professional	Improved consistency, interoperability, and discoverability of data, enabling machine-readable metadata and	Standards, Capacity	Confirmed

¹² <https://www.coretrustseal.org/>

¹³ <https://fairsharing.org/>

		supporting reuse across systems and communities		
Recommend appropriate open and non-proprietary file formats for a dataset to support long-term accessibility and interoperability	Researcher, Data Support Professional	Researchers are guided towards open, self-describing formats suitable for sharing and reuse; format recommendations are provided across the data publication workflow	Standards, Capacity	Confirmed
Validate and verify metadata quality against community best practices, recommendations and standards, producing records that support data reuse	Researcher, Institution	Harmonised, machine-interpretable metadata records that meet applicable community best practices and recommendations and support data reuse	Capacity, Standards, Quality and Provenance	Confirmed
Draft and review data documentation, including READMEs and data dictionaries, across the data publication workflow	Researcher, Data Support Professional	Clear documentation that enables others to understand and reuse the data	Capacity, Support	Confirmed
Map metadata across multiple schemas so the same dataset can be described in more than one standard, enabling deposit across different repositories without manual remapping	Researcher, Data Support Professional	Metadata is expressed in multiple schema formats simultaneously, reducing barriers to deposit across different repositories and improving interoperability without requiring the researcher to manually remap their record	Standards, Capacity	Confirmed
Enable consistent, reusable metadata across institutions and between different deployments of the Data Director, so data created in one system can be understood and used in another	Researcher, Data Support Professional	Existing metadata formats and schemas are reused; less data wrangling; easier interpretation of data across systems	Standards, Capacity	Confirmed

Detect and flag information in datasets that could identify individuals, either directly or indirectly, before publication	Researcher, Data Support Professional	Potential identifiers are flagged before the publication workflow proceeds; the researcher is alerted and guided on appropriate next steps; the tool does not attempt anonymisation or deidentification autonomously	Governance, Policy	Future Consideration
Check that variables use correct and consistent terminology to support future meta-analysis and data reuse	Researcher, Data Support Professional	Variable terminology is flagged where it may affect interoperability or reuse; researchers are guided towards community-agreed definitions	Standards, Capacity	Future Consideration
Generate or assist with Model Cards or Dataset Cards for machine learning research datasets	Researcher, Data Support Professional	A structured Model Card or Dataset Card is produced for machine learning datasets, improving transparency and reusability. Note: this use applies specifically to machine learning datasets	Capacity, Standards	Future Consideration
Create crosswalks between metadata schemas to support interoperability	Researcher	Datasets from different sources can be combined and compared, for example medical imaging data with geospatial data	Standards, Capacity	Future Consideration
Harmonise metadata schemas across repositories to improve consistency and FAIR alignment across institutional and regional systems	Data Support Professional	Interoperable repositories at institutional, regional and national level	Standards, Capacity	Future Consideration
Connect multiple Data Director deployments so datasets across institutional boundaries can be discovered and shared through a joined-up system, including data from third-party sources	Researcher, Data Support Professional	Data held across institutional boundaries and in external sources becomes discoverable and shareable through connected Data Director instances	Standards, Capacity	Future Consideration

3. PERSISTENT IDENTIFICATION				
Support the assignment of persistent identifiers (PIDs) to datasets at the point of depositing data for publication; users are informed if the selected repository assigns PIDs automatically	Researcher, Institution	Datasets are citable and discoverable via persistent identifiers (stable, long-lasting references that remain valid if a dataset moves location); impact metrics can be collected; PID assignment remains the responsibility of the repository	Support, Standards	Confirmed
Support PID assignment at variable or file level, not only at dataset level, to enable more granular citation and traceability	Researcher, Data Support Professional	Individual variables or files within a dataset are persistently identifiable and citable; granular provenance and credit attribution is enabled	Standards, Quality and Provenance	Future Consideration
Attach machine-readable licence and access policies to datasets so that permitted uses travel with the data and can be understood by both humans and automated systems	Researcher, Data Support Professional	Licence and access conditions are machine-readable and travel with the dataset as a persistent, FAIR-aligned policy object	Governance, Standards	Future Consideration
4. DATA MANAGEMENT PLANNING				
Check that research data publication matches what was committed to in the Data Management Plan (DMP), where one exists	Researcher	Research data is published in line with the original plan; any discrepancies are flagged for human review; the tool functions without a DMP where one is not available	Capacity, Policy	Confirmed
Draft, review and update DMPs, incorporating FAIR and CARE principle considerations; CARE principles relate to the governance of Indigenous data.	Researcher, Data Support Professional	A DMP that identifies FAIR and CARE gaps, meets funder and institutional requirements, and is kept up to date throughout the project	Capacity, Policy, Governance	Future Consideration

Update the DMP to reflect research data as published, where funders require this	Researcher, Data Support Professional	The DMP accurately reflects the final data publication; funder compliance is maintained post-publication	Capacity, Policy	Future Consideration
5. METADATA REMEDIATION AND IMPROVEMENT				
Recommend improved metadata for already-published datasets and supplementary materials, including preliminary and unpublished datasets, for researcher review and action	Researcher, Data Support Professional	Existing datasets become more discoverable, citable, and reusable; metadata is improved against FAIR principles and community best practices	Capacity, Standards	Confirmed
Use automated data processing to extract, reformat and improve existing metadata records, applying the same standards-selection process used for new data	Researcher, Data Support Professional	Existing metadata is cleaned and updated using the same standards-selection process as new data	Capacity, Standards	Future Consideration
6. FAIR ASSESSMENT				
Connect to recognised FAIR evaluation tools to assess how well produced metadata meets community standards, returning a score or maturity level to guide improvement	Researcher, Data Support Professional, Institution	A FAIR evaluation score or maturity level is returned by the connected evaluation tool; researchers can understand, track, and improve their metadata against recognised criteria	Standards, Quality and Provenance	Future Consideration
Evaluate how FAIR-aligned already-published datasets are using recognised tools such as FAIRassist , ¹⁴ and provide improvement recommendations	Researcher, Data Support Professional	The FAIRness of existing published outputs is assessed against community benchmarks; improvement recommendations are provided	Standards, Quality and Provenance	Future Consideration

¹⁴ <https://fairassist.org/>

7. PROVENANCE AND TRACKING				
Track the provenance chain of AI agent actions on data and metadata	Researcher, Data Support Professional, Institution	All agent actions are recorded and can be verified and validated	Quality and Provenance, Standards	Confirmed
Produce provenance records for existing datasets and enable them to be cited and connected to related outputs, so data sources are traceable and credit is attributed	Researcher, Data Support Professional	Data sources are traceable; credit is attributed to the original dataset	Quality and Provenance, Standards	Future Consideration
Assign persistent identifiers and provenance records to individual data variables so each carries full context including definition, units, and origin, supporting human and machine use	Researcher, Data Support Professional	Individual variables are identifiable, citable, and carry full contextual provenance including definition, units, properties, and relationships	Quality and Provenance, Standards	Future Consideration
Digitally sign large language model prompts and verified question and answer pairs to record ownership and provenance of AI-generated content	Researcher, Data Support Professional	AI-generated content is traceable to its source prompt; ownership is recorded with persistent identifiers; verified human-approved outputs are signed and auditable	Quality and Provenance	Future Consideration
8. DISCOVERY AND LINKING				
Link datasets to related research outputs such as journal article publications, people, grants and software at the point of publication	Researcher, Data Support Professional, Institution	Persistent links are established; richer reporting across research assets; more contextual discovery of outputs	Standards, Support	Confirmed
Discover related datasets relevant to a research project	Researcher	Relevant data is surfaced quickly without relying on general-purpose search engines; discovery does not affect existing metadata	Standards, Support	Confirmed

<p>Include data citations in related journal article publications and record the relationship between a dataset and its associated publication in the metadata</p>	<p>Researcher, Data Support Professional</p>	<p>Datasets are linked to related journal article publications using metadata properties such as IsReferencedBy and IsSupplementTo; where the article cannot be updated, the relationship is recorded in the dataset metadata</p>	<p>Standards, Support</p>	<p>Future Consideration</p>
<p>Record deposit in institutional systems such as Current Research Information Systems (CRIS)</p>	<p>Researcher, Data Support Professional, Institution</p>	<p>Deposit is recorded in the relevant institutional Current Research Information System (CRIS), supporting institutional reporting on research outputs and activities</p>	<p>Policy, Support</p>	<p>Future Consideration</p>
<p>Monitor access requests, downloads and citations for published datasets and provide researchers with metrics on data reuse</p>	<p>Researcher, Institution</p>	<p>Researchers receive regular updates on how their data is being accessed and reused; metrics support recognition of data outputs within the current reward system</p>	<p>Standards, Support</p>	<p>Future Consideration</p>
<p>Link datasets to protocol repositories to enhance data replicability and methodological transparency</p>	<p>Researcher, Data Support Professional</p>	<p>Datasets are connected to the protocols used to generate them; replicability is supported through persistent links to methodology</p>	<p>Quality and Provenance, Standards</p>	<p>Future Consideration</p>
<p>Find journal article publications and linked datasets on a given research topic</p>	<p>Researcher</p>	<p>A bibliography of relevant papers and associated datasets is produced</p>	<p>Standards, Support</p>	<p>Future Consideration</p>
<p>Build knowledge graphs, for example from nanopublications</p>	<p>Researcher</p>	<p>Knowledge is shared and managed in a structured, machine-readable form</p>	<p>Standards</p>	<p>Future Consideration</p>

9. PRIVACY AND COMPLIANCE				
Carry out privacy due diligence for research involving human subjects at the point of publication	Researcher, Data Support Professional	Consent documents and data use agreements are checked; identifiers are removed; ethics approval documentation is linked; sharing is confirmed as legally compliant	Governance, Policy	Future Consideration
Manage data access automatically in line with applicable data protection legislation, for example restricting access to sensitive datasets, subject to explicit human approval	Data Support Professional	Data access is managed in line with legal requirements; no access restriction is enacted without human confirmation	Governance, Policy	Future Consideration
Monitor compliance and recommend remediation where data does not meet national or institutional policies	Data Support Professional	Non-compliance is identified; fixes are suggested, such as changes to metadata, file formats, or updates to related repository records.	Governance, Policy	Future Consideration
10. INDIGENOUS DATA				
Ensure data reflects Indigenous context and aligns with CARE principles at the point of publication	Researcher, Data Support Professional	Data is handled in a way that is not biased; Indigenous data rights and protocols are respected	Governance, Standards	Future Consideration
11. AI AGENT GOVERNANCE				
Establish digital identities for AI agents and manage their rights and permissions	Researcher, Data Support Professional	Humans can create agents with clearly defined, responsible rights; agents can return access and permissions to humans; data transfer is policy-driven and restricted appropriately	Governance, Quality and Provenance	Future Consideration

12. COMMUNITY OVERVIEW

Provide an overview of data practices within a research community	Researcher, Data Support Professional	A dashboard showing the metadata standards and most-used repositories within a community	Standards, Support	Future Consideration
---	---------------------------------------	--	--------------------	-----------------------------

5. Blueprint Considerations

This Blueprint sets out an ambitious range of capabilities, and implementers will encounter tensions and contradictions between principles and features that must be navigated according to context. This section sets out the conditions that any implementation must satisfy before the architecture principles in Section 6 can be meaningfully put into practice. It addresses six areas: governance of the Data Director as an agentic AI system, covering agent behaviour, automation boundaries and human oversight; governance of the data on which it operates, covering legal, ethical and sovereignty requirements; how to resolve conflicts between requirements; operational readiness for deployment; and costs, sustainability and evaluation. Together these provide the foundation from which all other implementation decisions follow.

Both agentic AI governance and research data governance are large, complex, and rapidly evolving fields that extend well beyond the scope of this Blueprint. The considerations set out here are not intended to be exhaustive but to provide implementers with a principled starting point for trustworthy and sustainable deployment. Automation should proceed where governance conditions have been defined and met, not simply because a capability is technically available. Beyond these foundational constraints, implementers have the freedom to configure, prioritise, and extend the Data Director's capabilities in ways appropriate to their institutional, disciplinary, and jurisdictional context, and are encouraged to engage with relevant national frameworks, institutional policies, and emerging community standards as these fields develop.

5.1 Governance and Roles

- **Define accountable human roles and responsibility assignments before configuring the tool.** Who approves consequential decisions, who resolves policy exceptions, and who maintains configuration must be documented in governance agreements. Where multiple institutions collaborate, a dedicated inter-institutional framework agreement should be in place before deployment.
- **Establish a policy priority hierarchy before implementation.** Where institutional, funder, journal, and disciplinary requirements conflict, a documented resolution order must exist. The tool can surface conflicts but cannot resolve them without human-defined precedence rules.
- **Determine the level of human oversight required for each workflow before configuring automation.** Distinguish between actions requiring explicit human approval and those requiring ongoing monitoring. Where a DMP exists, discrepancies

with actual data publication should trigger human review rather than automatic blocking, since a DMP is a plan and not a binding contract.

5.2 Legal, Ethical and Sovereignty Prerequisites

- **Treat legal, ethical, and human rights compliance as prerequisites for any sharing action.** Legal and ethical permissibility must be confirmed before workflows are enabled. Relevant ethics approval body status, such as an Institutional Review Board (IRB), Research Ethics Committee (REC), or equivalent national body, must be verified before data publication proceeds, and embargo status must be confirmed by the responsible human before any sharing workflow is triggered.
- **Protect sensitive data from unauthorised access or ingestion at all stages of the workflow.** Sensitive, personal, or legally restricted data must not be ingested by the large language model component of the tool without explicit institutional authorisation and appropriate legal safeguards. Technical controls must be in place to enforce this.
- **Resolve data residency and cross-border transfer requirements at the architecture stage.** Implementers must document how AI processing, logs, temporary files, and backups will be managed within jurisdictional boundaries. Practical approaches to achieving this depend on institutional infrastructure.
- **Ensure data sovereignty and community authority inform openness decisions from the outset, not retrospectively.** Indigenous data rights and national sovereignty obligations are conditions that shape what can be shared, not constraints to be accommodated afterwards. CARE principles must be applied as design inputs and must not be treated as a uniform standard, given significant variation across domains and communities.

5.3 Resolving Conflicting Requirements

Implementations will encounter genuine conflicts between requirements that cannot be resolved by the tool alone. Several architectural principles set out in Section 6 (Table 2) pull in opposing directions, and implementers must be prepared to make deliberate, documented decisions about how to resolve them. The following requirements apply to all implementations:

- **Document how conflicts between principles have been resolved and make that reasoning visible to users and auditors.** Resolution decisions must be recorded and communicated rather than left implicit in configuration choices.
- **Open data must remain openly accessible regardless of the security or encryption measures applied to the underlying infrastructure.** Encryption at rest and in transit protects infrastructure and must not be used to restrict access to datasets that are designated as open.

5.4 AI Behaviour and Automation

- **Ensure that agents do not operate anonymously.** All agent actions must be executed on behalf of a named, identified human user, and both the agent identity and

the human it acts for must form part of the provenance record. Where an agent acts on behalf of an institution rather than an individual, a designated accountable human role must be defined and documented.

- **Require that all agent actions on data and metadata are recorded and traceable, consistent with functional requirement R10 in Section 7.1.** A recognised provenance standard such as [PROV-O](#)¹⁵ should be adopted to ensure records are structured, interoperable, and auditable.
- **Where a dataset arrives with existing provenance information, the Data Director must preserve and carry forward that record rather than replacing or discarding it;** new agent actions are appended to the existing provenance chain.
- **Enable automation only where governance conditions have been defined and met, not simply because a capability is technically available.** Consequential decisions must remain with accountable humans. Override rights for AI outputs must be role-based and governed, not open to all users.
- **Where personal identifiers or sensitive information are flagged or detected, the tool must pause the sharing process rather than acting autonomously.** The researcher must be directed to appropriate external anonymisation or deidentification workflows before the sharing process can proceed.
- **Define explicit agent behaviour for when no controlled vocabulary, ontology, or community standard exists.** The tool must state the absence openly rather than falling back silently to general-purpose alternatives. In domains where no prior FAIR work has taken place, outputs may be entirely generic; implementations must detect and record this condition, communicate it clearly to the user, and include it in the provenance chain.
- **Prevent poor-quality AI-generated metadata from reaching repositories.** Quality checks on AI outputs must be a core capability, not an optional enhancement. Generic, uncurated, or inaccurate outputs must be flagged for mandatory human review before submission. This is distinct from FAIR assessment, which evaluates metadata against recognised external standards.

5.5 Operational Readiness

- **Define rollback and recovery mechanisms that do not compromise traceability or auditability.** Recovery to a known good state must be designed alongside provenance tracking so that interpretability and auditability are maintained through failure and recovery events.
- **Assess whether sufficient technical, policy, and data management capacity exists to configure and maintain the tool before proceeding with deployment.** This includes the availability of staff who can manage technical configuration, maintain policy mappings as requirements change, and provide human oversight of AI-generated outputs.

¹⁵ <https://www.w3.org/TR/prov-o/>

- **Guard against over-reliance on the tool, which can erode in-house research data management expertise over time.** Deployment should be accompanied by investment in researcher and data professional skills, as noted in Section 3.1.3.
- **Calibrate reliability and uptime targets to the deployment context.** Commitments appropriate for national infrastructure are not realistic for locally hosted or resource-constrained deployments.

5.6 Costs, Sustainability and Evaluation

- **Make all implementation costs explicit and plan for the ongoing sustainability of the deployment from the outset.** Costs to document include deposit fees, PID registration costs, per-unit AI processing costs, staff time for configuration and ongoing maintenance, and the risk of scaling costs as usage grows. Institutions should model costs at projected scale before committing to deployment.
- **Ensure that implementation decisions do not introduce features, add-ons, or dependencies that make the tool inaccessible to lower-resource institutions.** Cost barriers and irreversible technical dependencies must be actively avoided, consistent with the affordability requirement set out in Section 7.2 (C17).
- **Establish a FAIRness baseline before deployment and define what metrics will demonstrate value over time.**
- **Ensure the tool meets recognised digital accessibility standards and is operable in resource-constrained environments.** Recommendations produced by the tool must be evaluated for implicit bias to ensure underrepresented communities and non-dominant metadata standards are not systematically excluded.

6. Architecture Principles

The following 14 principles govern all architectural decisions in this Blueprint and should be applied consistently across all layers and components (Table 2). They should be read alongside the Blueprint Considerations in Section 5, which set out the conditions that any implementation must satisfy before these principles can be meaningfully put into practice. The two are complementary: where the principles set direction and design intent, the considerations define the boundaries within which that direction must be pursued. Where principles pull in opposing directions, Section 5.3 sets out how to resolve them. Implementers may extend this list, but should not remove or contradict these principles without community agreement.

Table 2. The 14 architecture principles governing all design decisions in this Blueprint, showing the governing rule, what each principle requires, and what it means for implementers in practice.

ID	Principle	Statement	Implication
	<i>Name of the governing rule</i>	<i>What the principle requires</i>	<i>What this means for implementers in practice</i>
P1	FAIR and CARE	All Data Director capabilities must actively support the Findability, Accessibility, Interoperability and Reusability of research data as a primary design objective. Where data relates to Indigenous or community-held knowledge, CARE principles (Collective Benefit, Authority to Control, Responsibility, Ethics) must be applied alongside FAIR.	All workflows, metadata schemas, PID assignment and repository recommendations must be designed and evaluated against FAIR criteria, with each FAIR dimension explicitly supported in implementation from the outset. While connection to community-recognised evaluation tools and RDA-approved benchmarks is a future capability, the architecture must be designed to support this connection when it becomes available. Where CARE principles apply, particularly in contexts involving Indigenous data or community-held knowledge, they must be treated as context-specific rather than universal. Their application requires direct engagement with the relevant Indigenous and community stakeholders and cannot be reduced to a checklist assessment.
P2	Open First	Data and metadata should be openly available unless a justified reason to restrict exists, respecting the principle of open as possible, closed as necessary. Being FAIR-aligned or technically shareable does not make data lawful to share. Legal and regulatory compliance is a condition for all data processing, sharing, and reuse actions.	Default sharing must be open; restrictions must reference a legitimate legal basis (for example, health data, national security, or data protection law). Closed-by-default is not compliant. Where restricted, a non-sensitive metadata record must remain openly accessible. Legal compliance must be verified before any sharing action. Regulatory obligations must be institutionally configurable. Where P2 conflicts with P10 (Data Sovereignty) or applicable legal obligations, Section 5.3 applies.

<p>P3</p>	<p>Security by Design</p>	<p>Security must be built into every component of the Data Director from the outset, not added retrospectively. Sensitive data handling constraints must be embedded in the architecture. The secure path must always be the easiest path for the user; where security and ease of use conflict, security takes precedence.</p>	<p>Default configurations must enforce encryption, role-based access, and audit logging. Provenance and infrastructure access are restricted by default. Restricted data must not be produced or accessible without explicit human-enabled access controls. Opt-out of security controls is not permitted. User roles must be explicitly defined before access controls are specified.</p>
<p>P4</p>	<p>Human in the Loop</p>	<p>AI agents may draft and suggest, but no consequential action may proceed without explicit human approval. The Data Director may assist and recommend but must not decide, authorise, or take responsibility for consequential decisions. Human-in-the-loop (requiring explicit approval before an action proceeds) and human-on-the-loop (ongoing monitoring with override capability) are distinct modes of oversight that must be differentiated and documented per workflow.</p>	<p>Human-review triggers must be defined for all consequential actions. Every irreversible action requires explicit human confirmation. Implementers must document whether human-in-the-loop or human-on-the-loop applies to each workflow. The tool must never operate as an autonomous authority over governance or responsibility.</p>
<p>P5</p>	<p>Full Traceability</p>	<p>Every action taken by an agent and every artefact produced must be traceable to its inputs, agent version, scripts, paths, and human authorisations. Provenance must be recorded as a directed graph, enabling full traceability across processes including large language model prompts, with export and import functionality.</p>	<p>Provenance must be recorded using a recognised community standard such as PROV-O (see Section 5.4). Digital fingerprints and signatures must be assigned to all resources and actions. Where human roles are recorded, a recognised taxonomy such as CRediT¹⁶ must be used. Human identifiers must use persistent identifiers (PIDs) such as ORCID. Provenance access is restricted by default. Dataset versions must be individually traceable. Records must support compliance evidence on demand.</p>

¹⁶ <https://credit.niso.org/>

<p>P6</p>	<p>Open Standards</p>	<p>The system must consume and produce metadata using open, community-maintained standards and specifications rather than bespoke formats, to facilitate interoperability and reusability across systems, disciplines, and scientific domains. It must align with RDA Recommendations and Outputs.</p>	<p>New schema requirements must be sought in existing open standards. Agents must distinguish between general-purpose and community-specific standards when making recommendations. Crosswalks and semantic mappings between standards must be documented to enable cross-domain data transition. Agents should recommend domain-appropriate standards drawing on community resources such as GO FAIR.</p>
<p>P7</p>	<p>Open Source</p>	<p>Implementations of the Data Director should be developed and distributed as open-source software under a permissive licence such as MIT or Apache 2.0. Open source is strongly preferred because it enables community contribution, independent audit of AI behaviour, and global reuse without access barriers.</p>	<p>Source code must be published in a publicly accessible repository following recognised coding standards. Contributions must be accountable. Dependencies should favour open-source components. Proprietary components must be explicitly justified, documented, and must not create lock-in. The open-source preference reflects P6 (Open Standards) and the Blueprint's technology-agnostic framing: openness is the default; this preference does not prohibit proprietary implementations.</p>
<p>P8</p>	<p>Explainable AI</p>	<p>All AI-generated recommendations, metadata records, and decisions must be accompanied by a human-readable explanation of how the output was reached. Explanations must be accessible to non-technical users. All prompts used in agentic AI and sources for metadata schemas must be transparent and publicly available.</p>	<p>Every AI action must surface a contextualised summary with links back to sources, including a visible chain-of-thought reasoning window. Explanations must be exportable for review. Agents must provide feedback on what can be improved or is missing for the targeted repository. Black-box outputs are not acceptable.</p>
<p>P9</p>	<p>Locally Adaptable</p>	<p>The Data Director must support meaningful configuration at the institutional level to reflect local policies, funder and journal requirements, and disciplinary practices, including approved repositories and PID systems, without requiring changes to the core specification.</p>	<p>Configuration points must be exposed for metadata schemas, repository preferences, access controls, protection levels, retention, storage backend, draft and publish workflows, review access, and integration with existing tooling. Core and configurable components must be architecturally separated. Policy inconsistencies must be surfaced explicitly. Core FAIR and security requirements cannot be overridden.</p>

<p>P10</p>	<p>Data Sovereignty</p>	<p>The Data Director must respect national sovereignty and data residency requirements, ensuring that data remains within required jurisdictional boundaries. This applies particularly to nationally funded research and sensitive data contexts. Sovereignty requirements extend beyond the primary dataset to include AI processing, prompts, logs, temporary files, backups, embeddings, telemetry, and any derived artefacts generated throughout the workflow.</p>	<p>Implementations must support local and edge deployment to meet data residency mandates without dependence on foreign cloud infrastructure. Compliance with national sovereignty requirements must be configurable at the institutional level and documented.</p>
<p>P11</p>	<p>Deploy Anywhere</p>	<p>The Data Director must be deployable across edge, local, and cloud environments, ensuring researchers can access and use the system regardless of their available infrastructure, technical resources, or institutional support. This principle is in direct tension with P10 (Data Sovereignty) in contexts where cloud deployment may be available but is jurisdictionally prohibited; Section 5.3 requires implementers to document how this tension has been resolved.</p>	<p>Implementations must not require cloud infrastructure as a prerequisite. A minimum viable deployment must be defined and documented for resource-constrained environments. No single deployment environment should be privileged in the core architecture.</p>
<p>P12</p>	<p>Inclusive Access</p>	<p>The Data Director must be designed to include rather than exclude, across data sources, communities, users, and digital access needs. FAIR principles must be informed by CARE principles from the outset, ensuring that Indigenous and community voices shape how their data is represented and acted upon.</p>	<p>Implementations must assess data source readiness to identify exclusion risks. Digital accessibility standards must be met across all interfaces. Indigenous and community stakeholders must be included in governance and design. Autonomous agent data source selection must be governed to prevent unintended exclusion of relevant datasets or communities.</p>
<p>P13</p>	<p>Built to Last</p>	<p>The Data Director must be built on actively maintained, community-supported tools and architectures that do not create unsustainable operational burdens. External services, repositories, PID services, and vocabulary services represent risks outside the Data Director's control. When external dependencies are temporarily unavailable, the system must continue to function for unaffected workflow steps, alert the user clearly, and allow work to be</p>	<p>Dependencies must be evaluated for long-term viability before adoption. No single vendor or platform should be a single point of failure. Deprecated components must have documented migration paths. Workflows must queue or pause rather than fail silently when external services are unavailable.</p>

		saved and resumed without data loss.	
P14	Low Impact	The Data Director must minimise its environmental footprint. Edge data centres near end-users are the preferred default for compute tasks over centralised cloud infrastructure. The Data Director must not place unsustainable demands on community repository infrastructure.	Edge deployment is the default unless centralised infrastructure is demonstrably more efficient, which must be justified and documented. Repository recommendations must consider the capacity and funding status of target repositories to avoid contributing to submission overload. Where cloud deployment is adopted under P11, implementers should document how the environmental impact has been considered and why edge deployment was not viable. Energy footprint should be captured in provenance records at the component and workflow step level to support evaluation and optimisation over time.

7. Functional Specification

This section defines what the Data Director must do and how it must behave. It is structured in two parts.

Section 7.1 sets out the functional requirements: the specific capabilities the Data Director must provide, ordered to follow the workflow sequence established in the uses table in Section 4.1. Each requirement is assigned a priority level and mapped to known existing tools and resources in the Existing Solutions column. The Existing Solutions column is informational, not prescriptive. It identifies tools and resources the community has confirmed are relevant or compatible, to support implementers in meeting each requirement. It is not a complete list and does not constitute an endorsement of any specific product or platform; it reflects the state of community knowledge at the time of writing. Entries where no confirmed existing solution has been identified are marked accordingly; where a known tool is in active development or awaiting community evaluation, this is noted against the relevant entry.

Section 7.2 sets out the non-functional requirements: the quality standards and constraints the Data Director must meet across all its capabilities, regardless of which functional requirements are implemented.

7.1 Functional Requirements

The following requirements describe the specific capabilities and behaviours the Data Director must provide (Table 3). They are ordered broadly to follow the workflow sequence established in the uses table in Section 4.1, though some requirements span multiple workflow phases.



The key words **MUST**, **SHOULD**, and **MAY** in this document are to be interpreted as described in [BCP 14](#) (RFC 2119, RFC 8174)¹⁷ when, and only when, they appear in capitals as shown here. Each is qualified using one of three terms:

- **MUST**: mandatory; no compliant implementation without this.
- **SHOULD**: strongly recommended; departure requires documented justification.
- **MAY**: optional but explicitly permitted.

Where a requirement is marked **SHOULD**, implementing organisations operating in contexts with more demanding compliance or audit requirements should consider whether that requirement should be treated as **MUST** in their deployment. This applies in particular to R1 and R8, as noted in the relevant rows.

¹⁷ <https://www.rfc-editor.org/info/bcp14>

Table 3. Functional requirements for the Data Director, showing each capability, its description, its priority, and known existing solutions or resources.

ID	Requirement	Description	Priority	Existing Solutions
	<i>The specific capability or behaviour</i>	<i>What the requirement means in practice</i>	<i>MUST / SHOULD / MAY</i>	<i>Relevant or compatible known tools and resources</i>
R1	Recommend and select repositories meeting funder, journal, institutional and community requirements	Identify suitable repositories and flag conflicts between funder, journal and institutional requirements. Repository selection should be initiated before metadata generation, since the selected repository determines the required metadata schema, controlled vocabularies, and validation rules. Recommendations must distinguish commercial from open repositories, surface deposit fees, and support affiliation-based suggestions.	<p>SHOULD</p> <p>Where repository recommendation is a mandatory service requirement, implementers should treat this as MUST.</p> <p>Note: implementing organisations must assess whether C3 Compliance (a non-functional requirement set out in Section 7.2, Table 4) requires treating R1 as MUST.</p>	<p>FAIRsharing.¹⁸ A curated registry of community-defined standards, databases and policies across disciplines; includes repository records and supports filtering by subject and data type. A FAIRsharing MCP server¹⁹ is available for direct integration with the Data Director.</p> <p>re3data.²⁰ A registry of research data repositories with filtering by subject, access type and deposit conditions.</p> <p>Cornell Data Storage Finder.²¹ An open-source tool for localised repository and storage recommendations; may serve as a reference implementation for R1.</p>

¹⁸ <https://fairsharing.org/>

¹⁹ <https://mcp.fairsharing.org/>

²⁰ <https://www.re3data.org/>

²¹ <https://finder.research.cornell.edu/>

<p>R2</p>	<p>Generate FAIR-aligned and discipline-appropriate metadata using controlled vocabularies, ontologies and community standards including CARE principles</p>	<p>Generate a draft FAIR-aligned metadata record, drawing on discipline-specific vocabularies and ontologies where available. Produce metadata meeting named community standards including domain-specific fields and CARE considerations where applicable. Researcher reviews before submission. Where no controlled vocabulary exists for a domain, the tool must explicitly state this; it should additionally suggest candidate terms drawn from related datasets or adjacent community terminology, clearly flagged as informal or unvalidated.</p>	<p>MUST</p> <p><u>Ontology Lookup Service (OLS)</u>.²² An ontology repository providing controlled vocabularies across disciplines.</p> <p><u>NFDI Terminology Service</u>.²³ A terminology service providing controlled vocabularies for research communities.</p> <p>OntoPortal instances. Including <u>BioPortal</u>²⁴ and <u>AgroPortal</u>,²⁵ domain-specific ontology repositories.</p> <p><u>Global Indigenous Data Alliance (GIDA)</u>.²⁶ Community standards with explicit CARE guidance for Indigenous data.</p> <p><u>RDA kernel metadata recommendations</u>.²⁷ Define essential domain-agnostic metadata fields.</p> <p><u>FAIR2 metadata specification</u>.²⁸ A community-driven metadata specification built on MLCommons Croissant, relevant to this requirement.</p>
------------------	---	--	--

²² <https://www.ebi.ac.uk/ols/>

²³ <https://terminology.services.base4nfdi.de/provider>

²⁴ <https://bioportal.bioontology.org/>

²⁵ <https://agroportal.lirmm.fr/>

²⁶ <https://www.gida-global.org/>

²⁷ <https://doi.org/10.15497/RDA00031>

²⁸ fair2.ai/spec

<p>R3</p>	<p>Suggest appropriate vocabularies, ontologies and open data formats</p>	<p>Recommend suitable controlled vocabularies, ontologies and file formats, including field-level formats such as date/time standards, for a given dataset to improve consistency, interoperability and discoverability. Must support emerging standards and new mappings. Must distinguish between ontology alignment and controlled vocabulary concept linkage.</p>	<p>MUST</p>	<p>FAIRsharing. Covers terminologies, formats, identifier schemas and reporting guidelines across disciplines; see R1 for MCP integration note.</p>
<p>R4</p>	<p>Validate metadata quality and produce standardised records assessed against community best practices, recommendations and standards</p>	<p>Check metadata completeness and correctness against defined schemas. The tool should connect to repository-provided validation tools where available, and to community validation tools where repository-specific tools do not exist. It does not replicate full FAIRness evaluation. Metadata validation outputs must be traceable through the provenance chain (see R10).</p>	<p>MUST</p>	<p>RO-Crate validator,²⁹ Frictionless Data,³⁰ JSON Schema,³¹ XSD,³² SHACL,³³ Croissant Validator,³⁴ Schema validation tools applicable to this requirement.</p> <p>CEDAR Workbench,³⁵ OntoPortal/BioPortal,³⁶ Skosmos,³⁷ OLS, FOOPS!.³⁸ Ontology and vocabulary validation tools.</p> <p>fairassist.org.³⁹ A full list of community FAIR assessment tools, including those with OSTrails framework integration.</p>

²⁹ <https://rocrate-validator.readthedocs.io/en/stable/>

³⁰ <https://frictionlessdata.io/>

³¹ <https://json-schema.org/>

³² <https://hackolade.com/help/XSDXMLSchemaDefinition.html>

³³ <https://www.w3.org/TR/shacl/>

³⁴ <https://huggingface.co/spaces/luisoala/croissant-checker>

³⁵ https://doi.org/10.1007/978-3-319-68204-4_10

³⁶ <https://ontoportal.org/>

³⁷ <https://skosmos.org/>

³⁸ <https://catalogue.fair-impact.eu/resources/foops>

³⁹ <https://fairassist.org/>

				<p>MOMSI FAIRsharing collection.⁴⁰ Cross-domain benchmarking reference covering domain-specific and universal standards, confirmed useful for evaluating R3 and R4: fairsharing.org</p> <p>FAIRshake,⁴¹ F-UJI.⁴² Metadata quality evaluation tools for assessing AI-generated record quality against curation standards; relevant to C15.</p>
R5	Draft data documentation including READMEs and data dictionaries at publication	<p>Generate a draft of supporting documentation so that others can understand and reuse the data at the point of data publication. The tool can automatically draft elements it can infer from available data and metadata, including file structure, field names, and basic descriptions. Elements that cannot be inferred and require direct researcher input include variable definitions, units, missing value codes, study design details, and data collection procedures. A human-reviewed and approved version must be produced before submission.</p>	MUST	No confirmed existing solution identified at v0.1. Known solutions to be identified through community input.

⁴⁰ <https://fairsharing.org/5742>

⁴¹ <https://doi.org/10.1016/j.cels.2019.09.011>

⁴² <https://www.f-ujl.net/>

<p>R6</p>	<p>Enable interoperable metadata across institutions and Data Director instances</p>	<p>Reuse existing metadata formats and schemas to reduce data wrangling and enable consistent interpretation across systems and institutions. The tool must support mapping of metadata across multiple schemas so that the same dataset can be described in more than one standard, enabling deposit across different repositories without manual remapping by the researcher. Researchers must be able to view original metadata alongside mapped versions and crosswalks. Synchronisation between Data Director instances must be supported where multiple institutions collaborate.</p>	<p>MUST</p>	<p>FAIR Digital Object (FDO) framework.⁴³ Supports interoperability through typed, validated objects; multiple profiles enable metadata at domain-agnostic, domain-specific and application-specific levels simultaneously.</p> <p>Cross-Domain Interoperability Framework (CDIF).⁴⁴ Supports cross-domain metadata mapping.</p>
<p>R7</p>	<p>Support assignment of persistent identifiers (PIDs) at point of publication</p>	<p>The Data Director generates PID-ready metadata and informs users of repository PID capabilities. It does not interact directly with PID registration systems; PID assignment remains the responsibility of the repository. Users must be informed whether the selected repository assigns PIDs, including whether DOIs are minted automatically. Any costs associated with PID assignment must be surfaced to the user and must not be absorbed silently into the application; institutional configuration</p>	<p>MUST</p>	<p>DataCite.⁴⁵ Provides the primary membership model for institutional DOI minting for deposited datasets; the Data Director should surface whether a selected repository holds DataCite membership.</p> <p>FAIRsharing identifier schema catalogue.⁴⁶ Maintains curated records describing identifier schemas and their persistence, resolvability, and global uniqueness characteristics, enabling repositories and applications to assess</p>

⁴³ <https://fairdo.org/>

⁴⁴ <https://worldfair-project.eu/cross-domain-interoperability-framework/>

⁴⁵ <https://datacite.org/>

⁴⁶ https://fairsharing.org/search?fairsharingRegistry=Standard&recordType=identifier_schema&page=1

		must determine how such costs are handled.		the suitability of identifier systems against EOSC PID criteria . ⁴⁷ Records are linked to the repositories that implement each schema.
R8	Verify data sharing at publication against Data Management Plan (DMP) commitments	At the point of data publication, the Data Director should verify that data sharing aligns with the associated DMP where one exists. Discrepancies must be flagged for human review, not automatically enforced. The system must function without a DMP where one is not available. Embargo status must be confirmed by the responsible human before any sharing workflow proceeds, consistent with Section 5.2.	SHOULD Where DMP compliance is an audit requirement, implementers should treat this as MUST.	Data Stewardship Wizard (DSW) . ⁴⁸ Exposes DMP commitments in a computationally amenable form; confirmed compatible with the Data Director. Argos . ⁴⁹ DMP tool with machine-readable output; confirmed compatible. DAMAP . ⁵⁰ DMP tool with machine-readable output; confirmed compatible. RDA maDMP Working Group . ⁵¹ For a complete list of compatible tools, contact the RDA Common Application Programming Interface (API) for Machine Actionable Data Management Plans Working Group.

⁴⁷ <https://fairsharing.gitbook.io/fairsharing/record-sections-and-fields/general-information/globally-unique-persistent-and-resolvable-identifier-schemas>

⁴⁸ <https://ds-wizard.org/>

⁴⁹ <https://argos.openaire.eu/home>

⁵⁰ <https://damap.org/>

⁵¹ <https://www.rd-alliance.org/groups/common-application-programming-interface-api-for-machine-actionable-data-management-plans-madmps/activity/>

<p>R9</p>	<p>Retrofit FAIR-aligned metadata for already-published datasets and supplementary materials</p>	<p>Generate metadata improvement recommendations for previously published, preliminary, or unpublished datasets to support more discoverable, citable, and reusable outputs. All recommendations are advisory and subject to human review and action before any changes are made. Automated repository updates are not within scope; the researcher or Data Support Professional is responsible for submitting approved recommendations to the relevant repository. Where a dataset lacks a persistent identifier, users should be directed to R7.</p>	<p>SHOULD</p>	<p>FAIR Digital Object (FDO) framework. Enables retrofitting by managing metadata independently of repositories, allowing harmonised metadata across multiple repositories without repository operator involvement.</p> <p>COMET initiative.⁵² Supports community-curated enrichment of PID metadata for existing outputs, directly aligned with the retrospective scope of this requirement.</p>
<p>R10</p>	<p>Track provenance of AI agent actions on data and metadata</p>	<p>Record all actions taken by AI agents on data and metadata so that they can be verified, validated and audited.</p>	<p>MUST</p>	<p>PROV-O.⁵³ Identified as a candidate standard for provenance recording, with archiving option. See also Section 5.4.</p> <p>eScience 2025 provenance paper.⁵⁴ Community reference for PROV-O implementation in agentic AI research workflows.</p> <p>yProv.⁵⁵ Provenance tracking service and includes yProv4ML.</p>

⁵² <https://www.cometadata.org/>

⁵³ <https://www.w3.org/TR/prov-o/>

⁵⁴ doi.org/10.1109/eScience65000.2025.00093

⁵⁵ <https://github.com/HPCI-Lab/yProv>

<p>R11</p>	<p>Link datasets to related research outputs including journal article publications, people, grants and software</p>	<p>The Data Director generates metadata that enables persistent links to be established between datasets and related research outputs including journal article publications, grants, software, and people. Linking metadata is produced for researcher review; registration of relationships through PID and repository workflows remains the responsibility of the researcher and the relevant repository. Linking to ORCID and ROR should be supported as part of metadata creation.</p>	<p>SHOULD</p>	<p>ORCID.⁵⁶ Persistent identifier for researchers; identified as a relevant linking target.</p> <p>ROR (Research Organisation Registry).⁵⁷ Persistent identifier for research organisations; identified as a relevant linking target.</p>
<p>R12</p>	<p>Discover related datasets relevant to a research project</p>	<p>Surface relevant existing datasets quickly without relying on general-purpose search engines. Discovery must be optional and must not trigger any changes to existing metadata records.</p>	<p>MAY</p>	<p>DataCite MCP. In active development at the time of writing; monitor for availability. Identified as a relevant resource for dataset discovery.</p>

7.2 Non-functional Requirements

Non-functional requirements describe how well the system must perform: the quality standards and constraints it must meet across all capabilities, such as how fast, secure, available, and accessible it must be (Table 4).

⁵⁶ <https://orcid.org/>

⁵⁷ <https://ror.org/>

Table 4. Non-functional requirements for the Data Director, showing each quality standard or constraint, its category, its requirement, and the rationale for its inclusion.

ID	Category	Requirement <i>What the system must do or achieve</i>	Rationale <i>Why this standard or constraint is necessary</i>
C1	Security	Data in transit and at rest must be encrypted at the infrastructure layer, with restricted access as default. Encryption operates at the infrastructure layer and does not restrict access to open data. Access to metadata, datasets and agent actions must be role-based, token-verified and auditable. A recognised framework such as NIST should be adopted as the compliance baseline.	The Data Director handles sensitive research data across multiple institutions. Funders, journals and institutions require demonstrable protection against unauthorised access, and audit trails are essential where AI agents act on data. The scope of the security audit trail is infrastructure and access control; AI agent action audit trails are addressed in C13. Australian Cyber Security Centre . ⁵⁸ Guidance on careful adoption of agentic AI services, directly relevant to agentic AI security baseline requirements
C2	Privacy	The Data Director must apply data minimisation principles to all metadata outputs, ensuring that personally identifiable or sensitive information is not exposed in metadata records or documentation. The tool must prompt users to confirm that applicable ethics approvals are in place before any deposit workflow proceeds. Where potential identifiers are flagged or detected in data or metadata, the tool must pause the sharing workflow, alert the human reviewer, and provide guidance on appropriate next steps; it must not attempt anonymisation or deidentification autonomously. Active detection of personal identifiers within datasets is identified as a future capability in Section 4.1; in this iteration the tool's contribution is to ensure that governance steps around personal data are not bypassed. More advanced detection of indirect or contextual identifiers remains a Future Consideration.	The Data Director's role is to detect and flag privacy risks, not to resolve them. Anonymisation or deidentification and ethics compliance are institutional responsibilities that require human judgement and context-appropriate tooling; the tool's contribution is to ensure these steps are not bypassed.

⁵⁸ cyber.gov.au/business-government/secure-design/artificial-intelligence/careful-adoption-of-agentic-ai-services

C3	Compliance	The Data Director must support compliance with funder open data requirements, applicable data protection legislation, institutional data policies, journal data policies and relevant disciplinary community standards. Evidence of compliance must be available on demand across all applicable obligations.	Funder mandates, journal policies and institutional requirements increasingly require open, FAIR-aligned data as a condition of grant award or publication. The Data Director must enable researchers and institutions to demonstrate compliance across all applicable obligations without significant manual effort.
C4	Sovereignty	Data and metadata must be stored and processed in jurisdictions compliant with institutional and national sovereignty requirements. Cross-border data transfer must be policy-controlled and auditable. Physical storage location must be constrainable. Location-based access controls must be supported.	Institutions operating under national and regional data regulations require assurance that research data does not move outside permitted jurisdictions without appropriate controls, especially for sensitive, embargoed, commercially restricted or Indigenous datasets where sovereignty obligations may be legally or ethically binding.
C5	Interoperability	The system must support established metadata standards, expose open APIs and enable metadata exchange with major repositories and institutional systems without bespoke integration. All Data Director instances must expose a harmonised API for core features to enable inter-Director interactions.	The Data Director's value depends on working across heterogeneous repository and institutional environments globally. Proprietary formats or closed interfaces would undermine FAIR alignment and limit adoption across systems, regions and institutions worldwide.
C6	Accessibility	The Data Director must meet recognised international accessibility standards across all user interfaces, including WCAG 2.1 as a minimum baseline. Documentation and guidance must be available in accessible formats and comply with applicable national accessibility legislation.	The Data Director will be used by a diverse global research community. Accessibility compliance is a legal requirement in many jurisdictions and reflects a commitment to inclusive research practice.

C7	Reliability	The Data Director must achieve a minimum of 99.5% uptime during agreed service hours for service-based deployments. Implementing organisations deploying the Data Director as a local or on-demand tool should calibrate reliability targets to their deployment context and document any departure from this target with justification, consistent with Section 5.5. All metadata submission and deposit workflows must complete successfully or fail gracefully with a clear error state, providing justification and suggested next steps.	Researchers rely on the Data Director at critical points in the publication workflow. Unplanned downtime or silent failures at point of submission could delay publication and compromise funder compliance. Reliability and availability are related but distinct: C7 addresses successful completion of workflows and uptime targets; C9 addresses the hours and time zones across which the system must be accessible.
C8	Resiliency	The Data Director must recover automatically from component failures without data loss. All AI agent actions on metadata must be reversible, with rollback to the last verified metadata state supported. Any failure mid-process must not result in corrupt or partial metadata records that compromise data integrity. Where an action has been explicitly confirmed as irreversible by an authorised human under C13, recoverability to a prior state may not apply; this must be logged and disclosed to the user at the point of confirmation.	Metadata creation and transformation workflows involve multi-step agentic AI processes where partial completion is potentially more harmful than complete failure. Corrupt or incomplete metadata records could propagate errors across repositories and undermine researcher trust in AI-generated outputs.
C9	Availability	The system must be available 24/7 to support researchers across time zones and at submission deadlines. Planned maintenance windows must be communicated in advance and minimised.	Research publication and data deposit activity is not confined to working hours. Grant deadlines and journal publication schedules require the Data Director to be accessible at any time. 24/7 availability also enables international cooperation across institutions.
C10	Performance	Metadata generation, validation and repository recommendation must provide visible progress feedback. Response times must be appropriate to operation complexity. Bulk operations must be processed asynchronously with user notification on completion.	Researchers interact with the Data Director at critical points in the publication workflow. Performance must be predictable and transparent. Slow or unpredictable performance without feedback will reduce adoption and create bottlenecks at publication.

<p>C11</p>	<p>Scalability</p>	<p>The Data Director must support growing volumes of datasets, metadata records and concurrent users without performance degradation, and must scale horizontally to meet regional or national demand. Infrastructure must be defined and managed as code to ensure consistent, reproducible and auditable deployment across institutional contexts. Infrastructure as Code approaches support version control of deployment configurations, reduce configuration drift between instances, and enable sovereignty-compliant deployment through jurisdiction-specific configuration profiles.</p>	<p>The Data Director is intended for adoption across multiple institutions and research communities. Usage will grow significantly as FAIR adoption increases and more repositories are onboarded. Individual Data Director instances will have capacity limits; implementations must be designed to support federation and growth across the wider ecosystem. Infrastructure as Code is particularly important in a multi-institution ecosystem where deployment consistency directly affects interoperability, reliability and data sovereignty compliance. The cost implications of horizontal scaling must be considered alongside C17 (Affordability).</p>
<p>C12</p>	<p>Data Governance</p>	<p>The Data Director must enforce institutional data governance policies including retention schedules, access controls, provenance tracking and audit logging. Governance policies must be machine-readable and AI-ready, drawing on frameworks such as ODRL. Compliance with well-recognised international data governance standards must be evidenced.</p>	<p>The Data Director operates across institutional boundaries. Each institution must apply and evidence its own governance policies while interoperating with shared community standards. A dedicated governance framework agreement between participating institutions may be required. The scope of the governance audit trail is institutional policy compliance; security audit trails are addressed in C1.</p>
<p>C13</p>	<p>AI Governance</p>	<p>All AI agent actions must be logged, attributable to a defined agent identity and subject to human oversight. Irreversible actions including deletion and autocompletion must be blocked by default without explicit human confirmation. The Data Director must comply with recognised international AI governance frameworks and standards. See also Section 6, Principle P4, for the distinction between human-in-the-loop and human-on-the-loop oversight modes.</p>	<p>The Data Director uses agentic AI to create, transform and validate metadata. Responsible AI use requires transparent, auditable and controllable agent behaviour in line with emerging national and international AI governance frameworks. Governance must be principle-driven across both data and AI dimensions. See also C14 (Explainability), which addresses human understanding and overrides rights for AI-generated outputs.</p>

C14	Explainability	The system must provide human-readable and understandable explanations for all AI-generated metadata, recommendations and decisions. Designated human-in-the-loop roles must be able to challenge and override AI outputs within a defined governance layer. A log of agentic AI reasoning steps and decision points must be maintained alongside metadata records for transparency and improvement.	Researchers and data support professionals need to trust and verify AI-generated outputs before submission. Opaque recommendations reduce confidence and increase the risk of errors in published metadata. Override rights must be role-based and governed rather than open to all users. See also C13 (AI Governance), which addresses logging, attribution and oversight of AI agent actions.
C15	Data and Metadata Quality	The Data Director must perform quality checks on its own metadata outputs and profile the quality of researcher-supplied input data. AI-generated metadata must meet the same curation standards as manually produced records. The system must evaluate curation readiness before submission. Low-quality or generic outputs must be flagged before submission with human review mandatory.	Low-quality AI-generated content is increasing across the internet. Metadata is required to be highly curated and precise. Preventing the propagation of generic, inaccurate, or uncurated AI-generated metadata into repositories is a pressing concern. Quality profiling of both AI-generated and researcher-supplied data ensures the system adds value rather than propagating poor-quality inputs.
C16	Measurability	The Data Director must collect usage data to measure the impact and usefulness of the tool, provide data for cost and benefit analysis, and track researcher uptake. Metrics must inform future development and enable institutions to justify involvement and subscription.	The effectiveness and impact of the Data Director must be demonstrable to justify institutional investment and inform future iterations. Uptake of digital research infrastructure requires evidence of value added for researchers across participating institutions and communities. Usage data collection must be designed in compliance with C2 (Privacy) and C4 (Sovereignty).

<p>C17</p>	<p>Affordability</p>	<p>The Data Director must be modular and configurable to support different funding levels and institutional contexts. Agents must not duplicate work, reusing outputs already created or discovered to reduce AI processing costs. Configurable tiers must ensure the system remains accessible to institutions with limited resources and without dedicated data support professionals.</p>	<p>The Data Director must not increase the digital divide. Institutions in lower-resource contexts, including those without human data support, must be able to access and sustain the tool. The balance between costs and advantages must be explicitly considered in architectural decisions. The scaling demands of C11 must be weighed against affordability to ensure infrastructure costs do not exclude lower-resource institutions.</p>
-------------------	-----------------------------	--	---

8. Process Flow

This section presents the Data Director workflow as a six-phase process flow. The process flow translates the functional requirements set out in Section 7.1 into a practical, sequential representation of how the Data Director supports researchers and data support professionals across the research data publication lifecycle. It is intended to help implementers understand how the tool's capabilities fit together as a coherent workflow, and to give researchers and data support professionals a clear picture of when and how the tool becomes relevant at each stage of their work.

The process flow was developed iteratively by Working Group members across two facilitated sessions using [Miro](#),⁵⁹ a collaborative online whiteboard, with participants contributing directly to the design in real time by adding steps, questioning sequencing, and refining individual process boxes. The guidance notes visible throughout the diagrams are drawn directly from Working Group contributions and record the reasoning behind specific design decisions. The workflow is therefore not a post-hoc illustration of the functional requirements but a community-designed artefact, developed in parallel with the specification and shaped by the practical knowledge of researchers, data support professionals and implementers across the Working Group.

The workflow is not prescriptive about technology or platform. It describes what happens, in what order, and where human decisions are required, without specifying how any individual implementation should deliver each step. Implementers should use it as a reference for designing workflows and configuring the tool within their own institutional context, guided by the implementation profiles in Section 11.

⁵⁹ https://miro.com/app/board/uXjVHXDM1lg=?share_link_id=139484060485

8.1 How to read the process flow

The workflow is structured in six phases, presented across six figures. The phases follow the sequence in which research data publication work typically unfolds, setting up before data collection; confirming governance and legal constraints; selecting a repository; preparing data and creating metadata; depositing and publishing; and managing outputs after publication. Together, the six phases cover the full journey from project planning to post-publication data reuse.

Researchers do not always enter the process at the beginning. Three entry points are provided to reflect this. Entry 1, shown in Figure 1, is for researchers starting a new project before data collection begins. Entry 2, shown in Figure 2, is for researchers whose data has already been collected and who are beginning to prepare it for sharing. Entry 3, shown in Figure 3, is for researchers who are ready to deposit data into a repository. All three paths converge and complete the same final phases.

Each figure uses a consistent set of visual conventions:

- Process steps are shown as plain rectangular boxes representing actions or tasks.
- Decision points are shown as diamond shapes with dashed orange borders, each presenting a yes/no question that branches the flow.
- Entry points are shown as rounded green ovals marking where a researcher joins the workflow.
- Phase headers are shown as orange-bordered bars that mark the transition between workflow phases.
- Functional requirement boxes, shown with an orange left border inside a process step, identify which requirement from Section 7.1 is being applied at that point.
- Guidance note boxes, shown with a blue left border, record session feedback and clarifying notes from the Working Group.
- Optional steps are shown with a dashed border.
- Arrows indicate the direction of flow throughout.

One requirement runs across all six phases without exception. R10, track provenance of AI agent actions on data and metadata, is shown as a red banner at the top of every figure. It applies to every action the Data Director takes throughout the entire workflow. This reflects the Blueprint's position that provenance recording is not a phase-specific capability but a foundational condition for trustworthy agentic AI operation. Every agent action must be recorded, and those records must be available for verification, validation and audit at any point.

8.2 Six-Phase Process Flow for Research Data Publication

The process flow described here represents an example pathway through the research data publication workflow as experienced by the researcher. It shows the stages a researcher moves through when preparing and depositing research data and identifies the points at which the Data Director can interact with and provide support to that process. The workflow belongs to the researcher; the Data Director's role is to assist, recommend, and flag, not to direct or

control. In practice, the sequence and emphasis of phases may vary depending on researcher context, institutional requirements, data sensitivity, and whether the dataset is being prepared for first-time publication or retrospective improvement.

Several design decisions embedded in the process flow are worth drawing attention to explicitly.

Repository selection, shown in Phase 3 (Figure 3), should be initiated before metadata generation begins, not only before deposit. The selected repository determines the required metadata schema, controlled vocabularies and validation rules that will govern all subsequent metadata work. Implementers should ensure that the tool prompts users to confirm or select a repository before the metadata creation workflow in Phase 4 is triggered (Figure 4).

Governance and constraints checks, shown in Phase 2 (Figure 2), precede repository selection and should inform all subsequent decisions. Sensitive data handling, ethical and legal constraints, embargo periods and DMP verification are all addressed at this stage so that the decisions made here can constrain and guide the steps that follow. This ordering reflects the Blueprint's position in Section 5.2 that legal and ethical permissibility must be confirmed before any sharing workflow is enabled.

Metadata quality validation, shown in Phase 5 as R4 (Figure 5), is placed before submission rather than after. This reflects Working Group feedback that quality issues should be identified and resolved while the researcher can still act on them, not discovered after data has been published.

The post-publication phase, shown in Phase 6 (Figure 6), is not the end of the process but a continuing loop. The Data Director may support embargo monitoring after deposit. It can also support retrospective metadata improvement for already-published datasets (R9) and the linking of datasets to related research outputs (R11). Implementers should design their deployments to remain active and useful to researchers beyond the point of deposit.

Phase 1: Pre-collection Setup

Data Publication Workflow | Entry 1: Starting a new research or data collection project

R10 (Horizontal Requirement — all phases): Track provenance of all AI agent actions on data and metadata. Record all actions so they can be verified, validated and audited.

LEGEND

- Process step**
An action or task for the researcher
- Decision point**
Yes/No — branches the flow
- Entry point**
Where a researcher joins
- Macro phase header**
Marks each workflow phase
- Functional requirement**
Agent capability (R1–R12)
- Guidance note**
Source labelled in each box
- Optional step**
Dashed border
- ▼ **Arrow**
Direction of flow

Entry 1: Starting a new research or data collection project

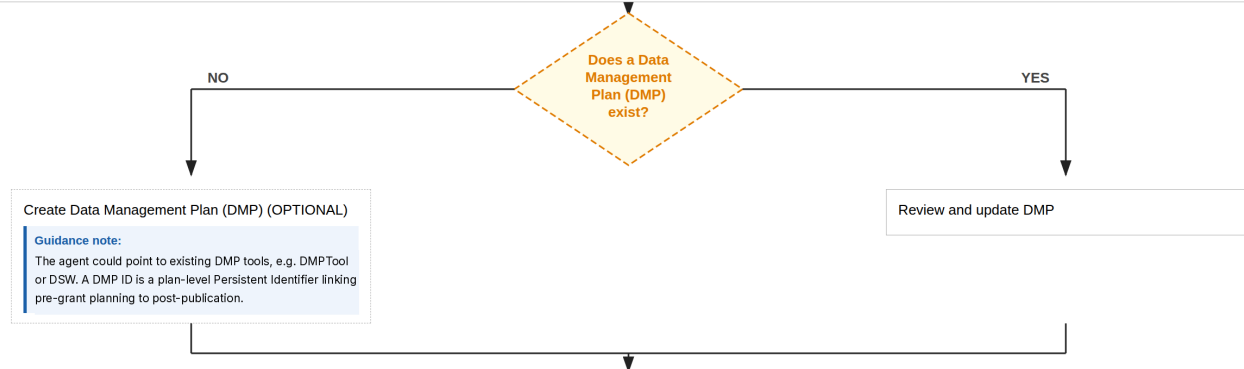
PHASE 1: PRE-COLLECTION SETUP (before data collection)

Draft data dictionary and README file

R5:
Draft data documentation including READMEs and data dictionaries at publication

Select controlled vocabularies and ontologies (e.g. via FAIRsharing)

R3:
Suggest appropriate vocabularies, ontologies and open data formats



→ Entry 2 (Data already collected) joins at Phase 2 (see next page)

Figure 1. Phase 1: Pre-collection Setup. The workflow entry point for researchers starting a new project before data collection begins. Shows the initial drafting of data documentation (R5), selection of controlled vocabularies and ontologies (R3), and a decision step to create a new Data Management Plan or review an existing one.

Phase 2: Governance and Constraints of Research Data

Data Publication Workflow | Entry 2: Data already collected joins here

R10 (Horizontal Requirement — all phases): Track provenance of all AI agent actions on data and metadata. Record all actions so they can be verified, validated and audited.

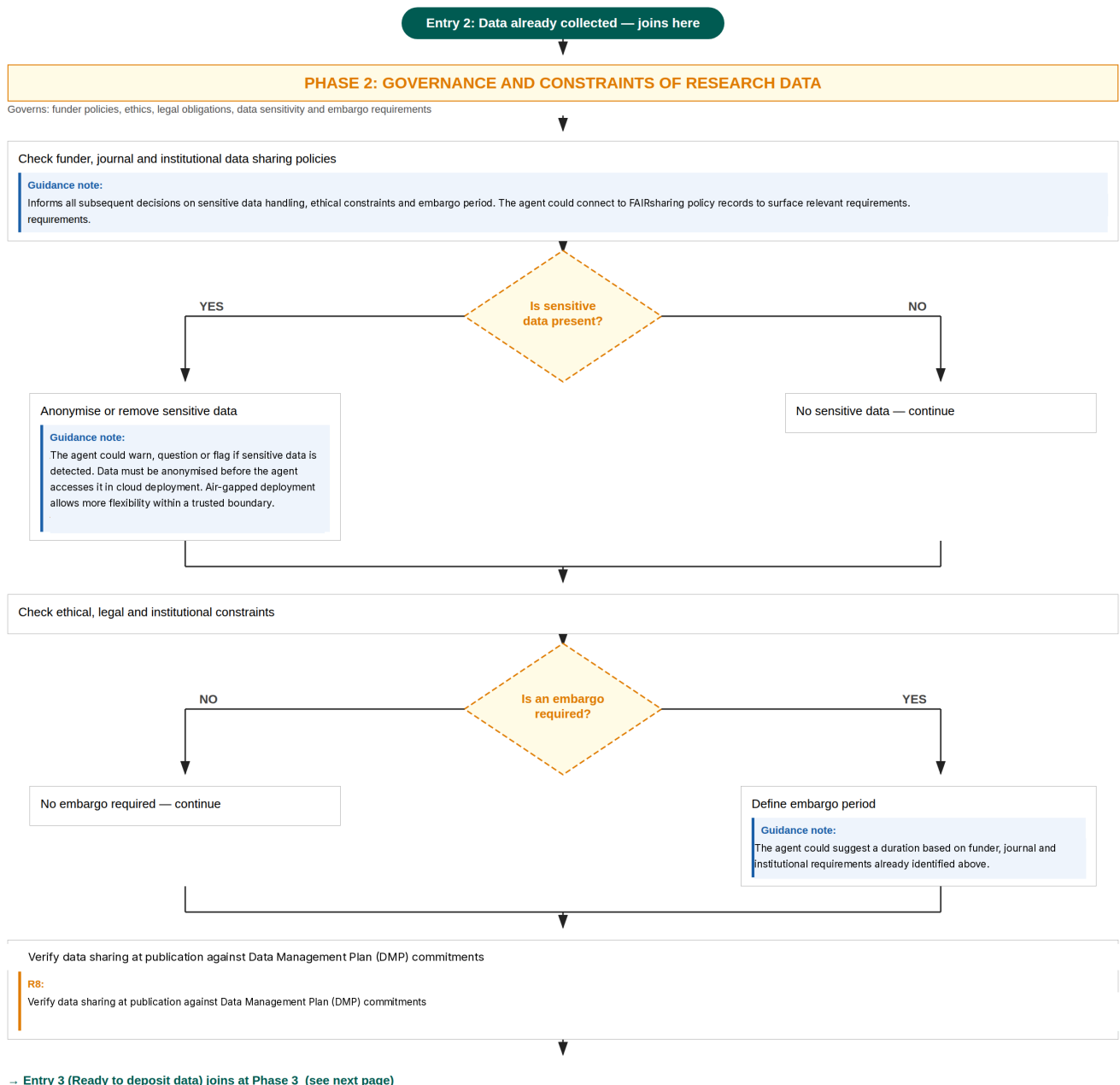


Figure 2. Phase 2: Governance and Constraints of Research Data. Covers the governance checks that must be completed before repository selection or data deposit proceeds. Includes checking funder, journal and institutional data sharing policies; identifying and handling sensitive data; confirming whether an embargo is required; and verifying data sharing commitments against the Data Management Plan (R8).

Phase 3: Repository Selection

Data Publication Workflow | Entry 3: Ready to deposit data joins here

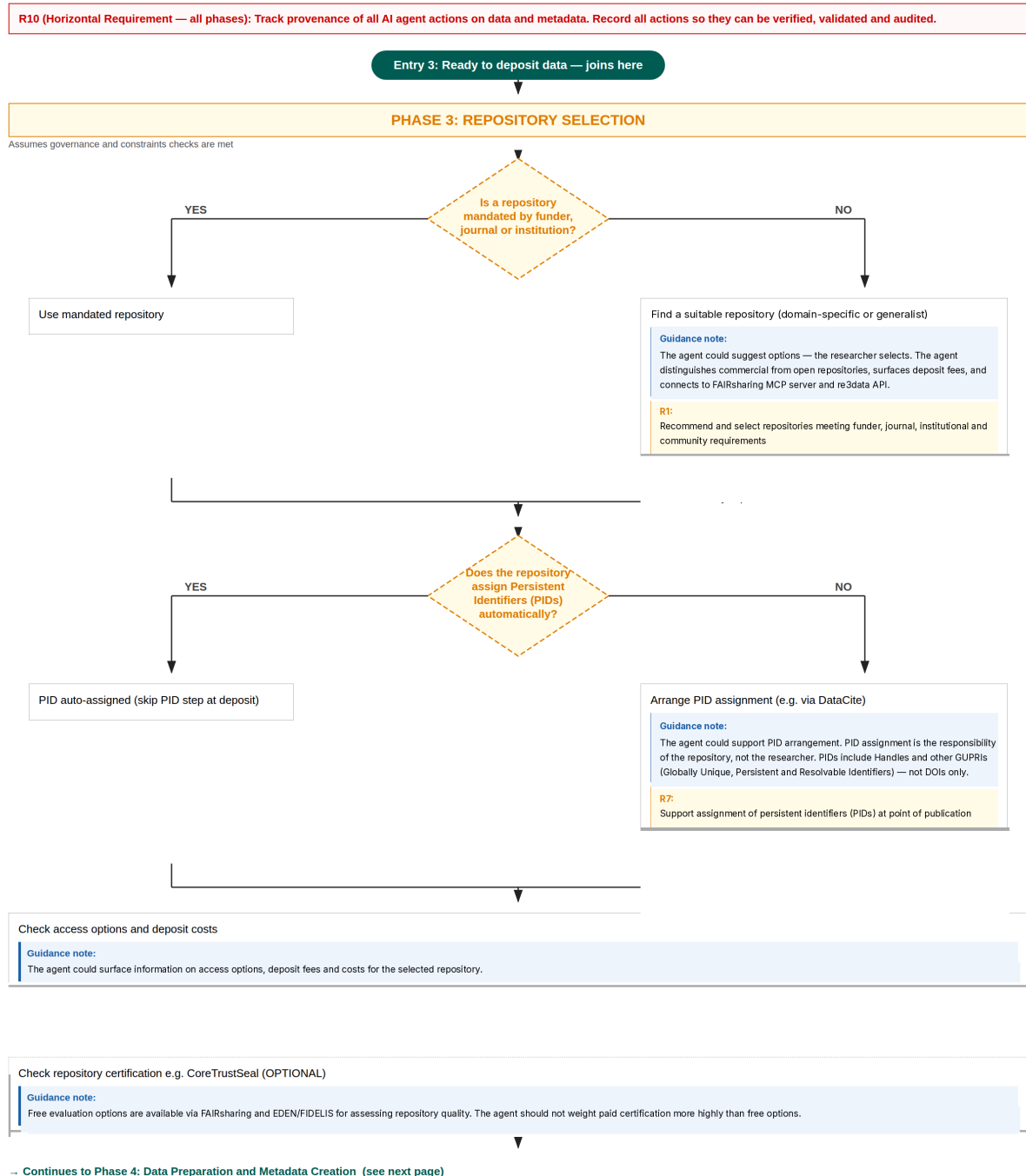


Figure 3. Phase 3: Repository Selection. Covers the selection and confirmation of a suitable data repository. Includes checking whether the funder, journal or institution mandates a repository; recommending an appropriate repository where one is not mandated (R1); confirming whether persistent identifiers (PIDs) are assigned automatically by the repository or arranging assignment where needed (R7); checking access options and deposit costs; and optionally verifying repository certification.

Phase 4: Data Preparation and Metadata Creation

Data Publication Workflow | Data preparation and metadata creation can run in parallel with each other

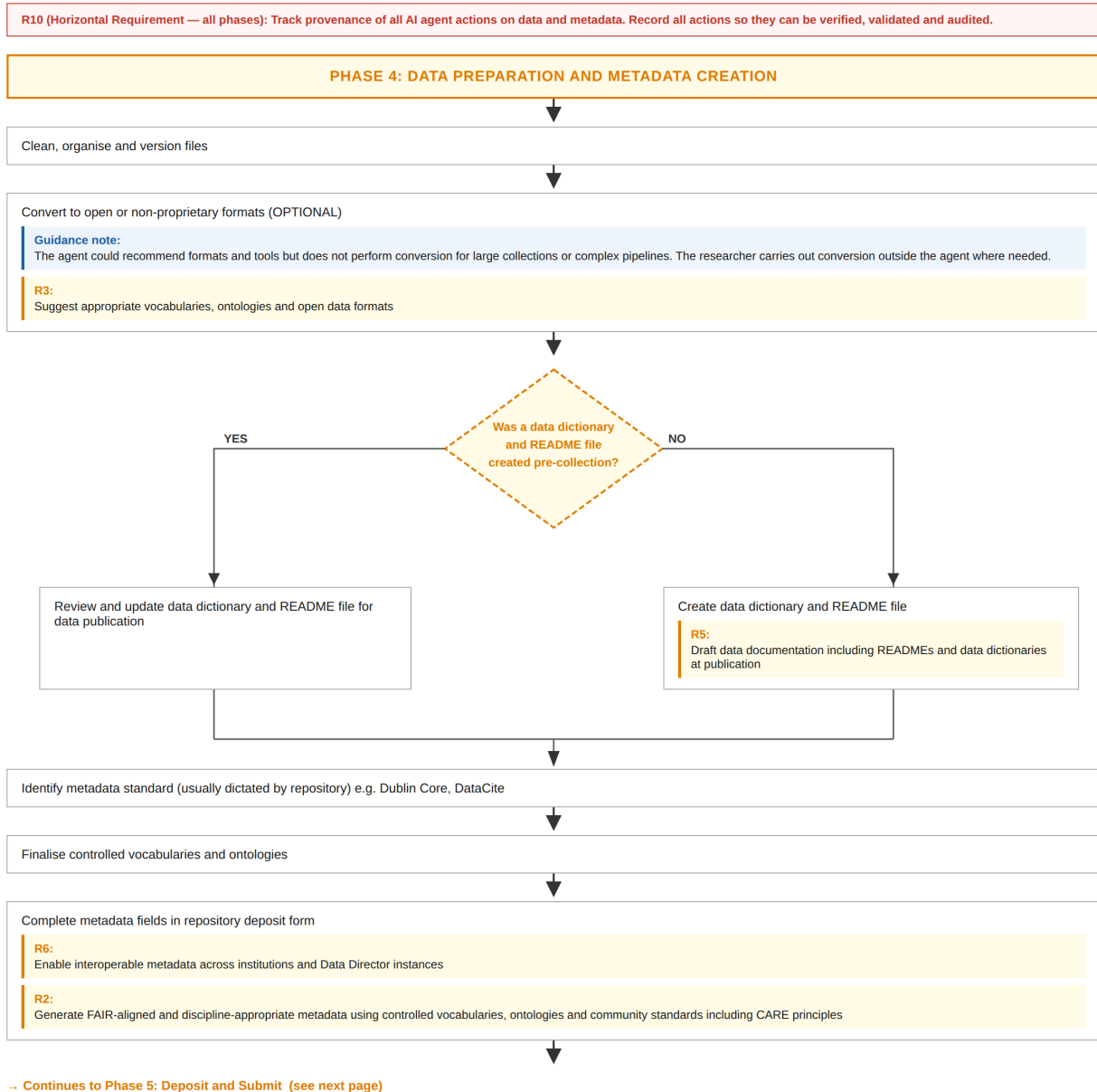


Figure 4. Phase 4: Data Preparation and Metadata Creation. Covers the preparation of data files and the creation of metadata records. Includes file cleaning and versioning; open format recommendations (R3); creation or review of data documentation, including READMEs and data dictionaries (R5); identification of the appropriate metadata standard; and completion of metadata fields using FAIR-aligned, discipline-appropriate vocabularies and ontologies (R2, R6).

Phase 5: Deposit and Submit

Data Publication Workflow

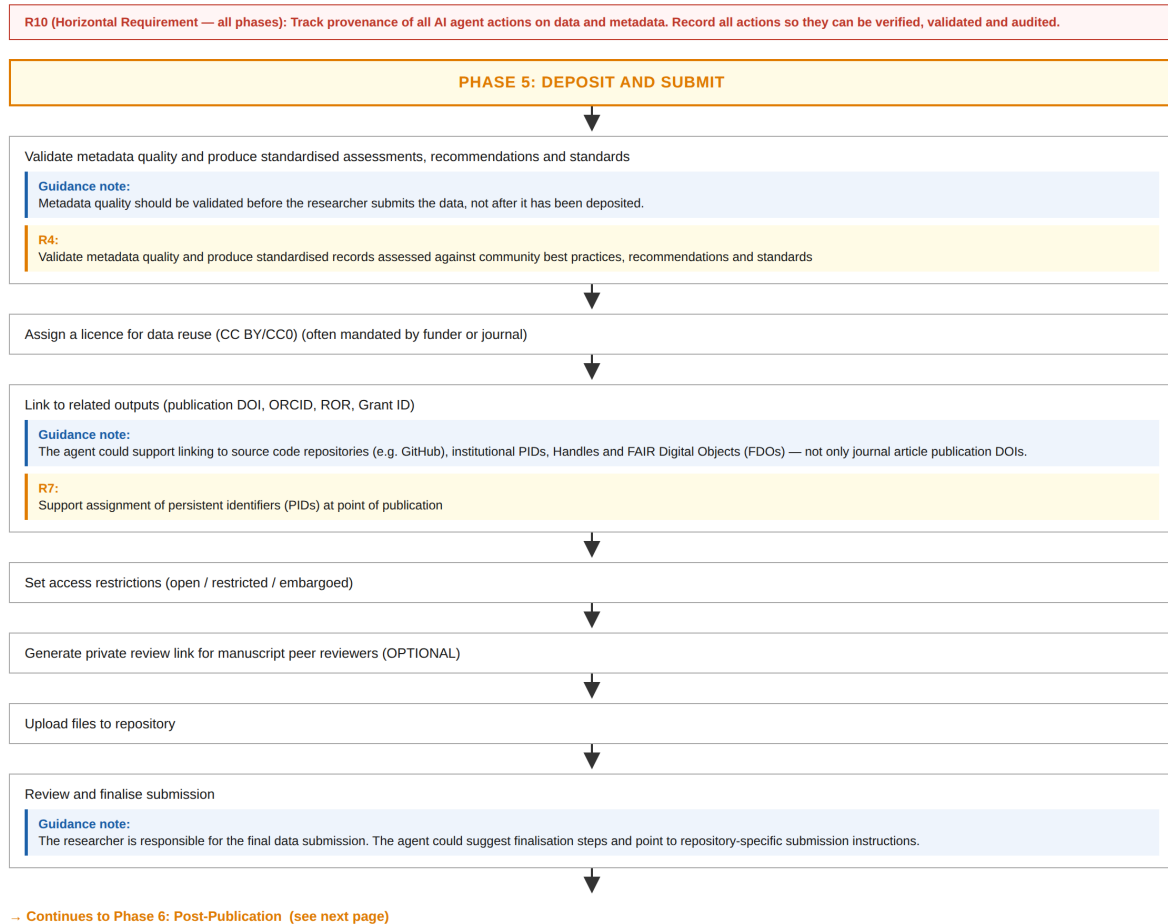


Figure 5. Phase 5: Deposit and Submit. Covers the final steps before a dataset is submitted to a repository. Includes metadata quality validation (R4); licence assignment; linking to related research outputs including publication DOIs, ORCID, ROR and Grant IDs (R7); setting access restrictions; generating a private peer review link (optional); uploading files; and researcher review and finalisation of the submission.



Phase 6: Post-Publication

Data Publication Workflow | Closes the loop back to planning

R10 (Horizontal Requirement — all phases): Track provenance of all AI agent actions on data and metadata. Record all actions so they can be verified, validated and audited.

PHASE 6: POST-PUBLICATION (closes the loop back to planning)

Note: 'publication' in this phase refers specifically to journal article publication unless otherwise stated.

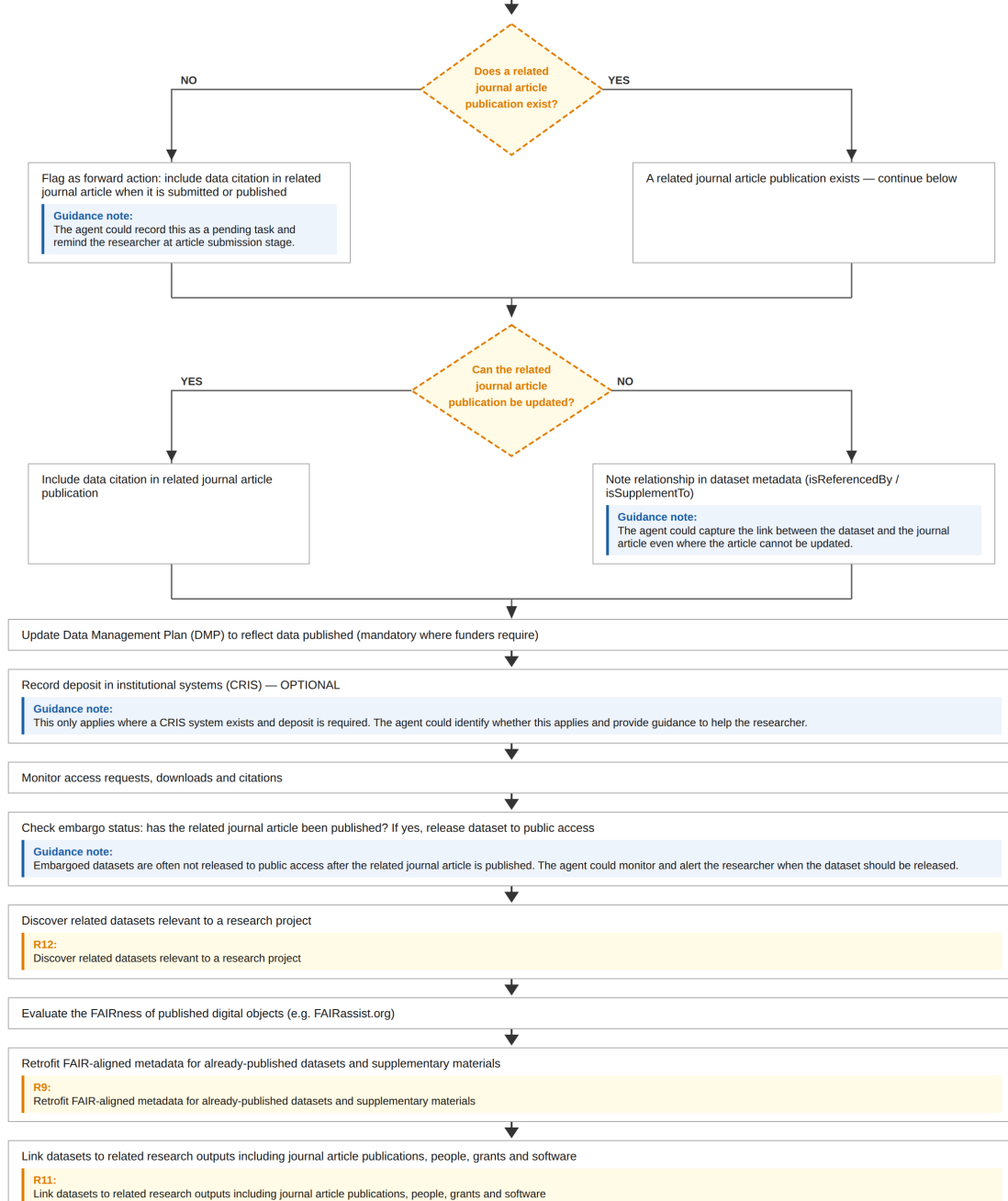


Figure 6. Phase 6: Post-Publication. Covers the actions that follow data deposit and closes the loop back to planning. Includes updating the Data Management Plan; recording deposit in institutional systems where required; monitoring access, downloads and citations; checking embargo release; evaluating the FAIRness of published digital objects and discovering related datasets (R12); retrofitting FAIR-aligned metadata for existing datasets (R9); and linking datasets to related research outputs (R11).

9. Reference Architecture

The architecture Blueprint provides a vendor-agnostic, technology-neutral map of the Data Director system components and how they interact. The Blueprint is also agnostic of the implementation scenario, such as cloud or hybrid cloud.

9.1 Architecture Overview

The diagram depicted below (Figure 7) illustrates the core capabilities required to support the data director implementation based on the principles and requirements identified during the working group discovery sessions.

Assumptions:

- Language support will be dependent on the chosen cloud/technology provider and requirements during implementation.
- High availability and disaster recovery requirements are dependent on the scenario and service level agreement (SLA) requirements during implementation.
- The architecture Blueprint is based on a greenfield implementation (a new implementation built from scratch) and does not account for all possible existing capabilities. It is recognised, however, that research organisations may already have services in place, including identity management and operational monitoring.

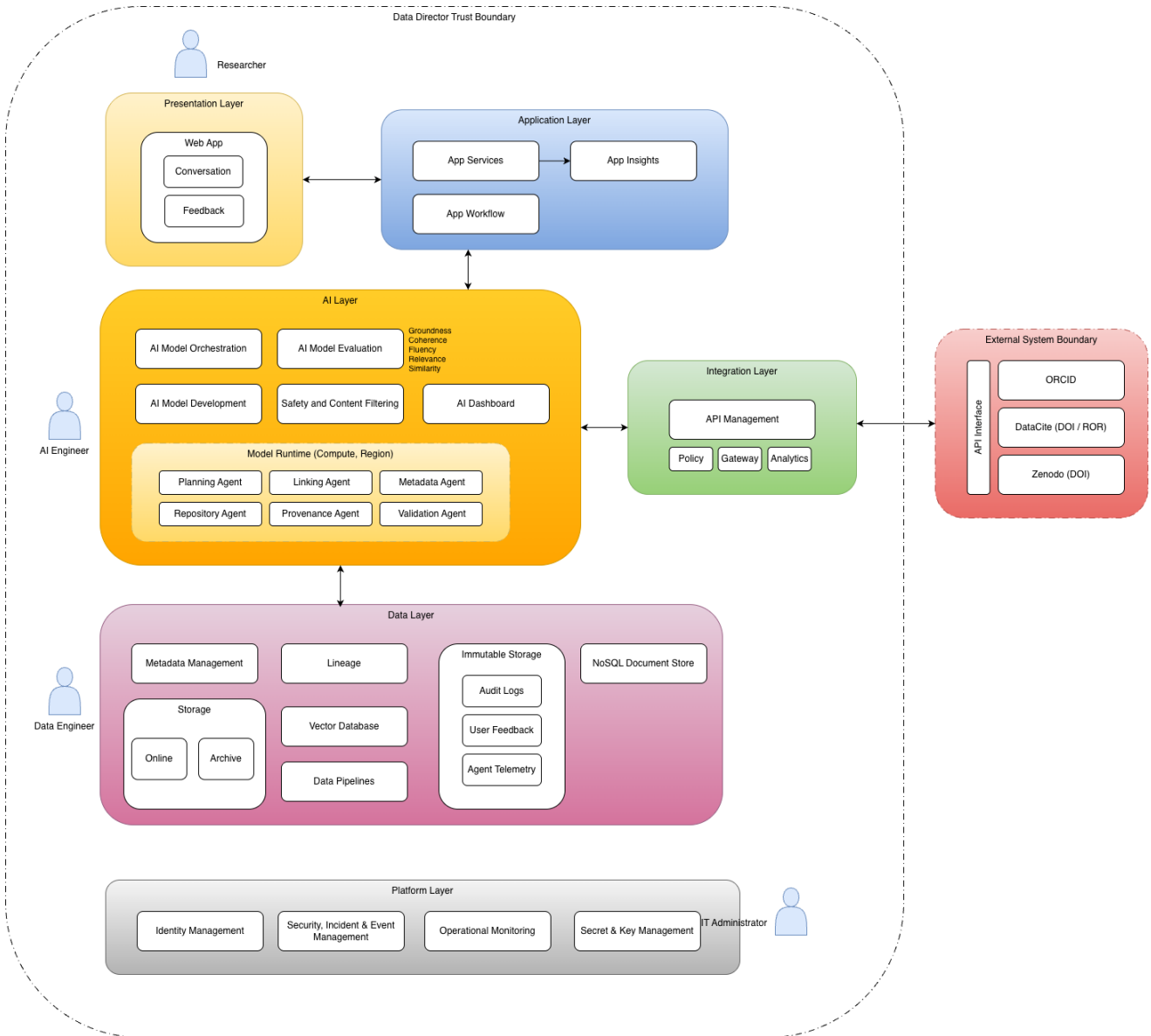


Figure 7. The Data Director conceptual architecture is a multi-layered AI platform that enables researchers to interact with intelligent agents for metadata, provenance, and repository management, underpinned by secure data storage, API integration with external scholarly systems, and governed platform infrastructure.

9.2 Presentation Layer

The presentation layer covers the interfaces through which users interact with the Data Director (Table 5). It must be accessible, responsive and clearly communicate AI-generated content separately from verified data.

Table 5. Presentation layer component, with its description and associated requirement, through which users interact with the Data Director

Component	Description/Requirement
Web app providing a chat interface and feedback mechanism	Supporting a chat-based interface which also allows the end-user to provide feedback on AI generated actions and output. Important note: the end user experience is subject to requirements specific to the research institution and implementation scenario and scope.

9.3 Application Layer

The application layer contains the core business logic of the Data Director: orchestrating user requests, managing workflows and coordinating between the AI and data layers (Table 6).

Table 6. Application layer components, with their descriptions and associated requirements, that orchestrate the Data Director's core business logic.

Component	Description/Requirement
App Services	Application Services which provide the compute and elastic scale to serve the end-user experience and interaction with underlying AI services.
App Insights	Application observability services which capture telemetry (usage, adoption, success, failure, audit information) which is useful for insights and troubleshooting.
App Workflow	Application-layer capability that orchestrates and manages end-to-end business processes across systems, users, and services in a reliable, stateful, and auditable way.

9.4 AI Layer

The AI layer contains all artificial intelligence and machine learning components (Table 7). This Blueprint does not prescribe specific models or vendors; it defines what the AI layer must do and how it must behave.

Table 7. AI layer components, with their descriptions and associated requirements, that provide the artificial intelligence and machine learning capabilities and agents underpinning the Data Director.

Component	Description/Requirement
AI Model Orchestration	AI capability responsible for coordinating, managing, and optimising the use of the AI agents and services to deliver reliable, consistent, and context-aware outcomes

AI Model Evaluation	AI capability responsible for systematically measuring, validating, and continuously assessing the performance, safety, and business effectiveness of AI models and agent behaviours
AI Model Development	AI capability which provides the end-to-end development and data science environment and processes for creating AI solutions, spanning data preparation, model selection, prompt and agent design, training or fine-tuning, and packaging models for deployment
Safety and Content Filtering	AI capability which enforces organisational, ethical, responsible and regulatory guardrails across all AI agent inputs and outputs to ensure that interactions remain safe, compliant, and appropriate for the intended use case
AI Dashboard	Provides monitoring, and operational insight into the performance, usage, cost, quality, and safety of AI models, agents, and workflows across an organisation
Model Runtime	Core application-layer capability responsible for executing AI models in production and managing the full lifecycle of inference requests in a reliable, scalable, and secure way
Planning Agent	Verifies DMP alignment, suggest open formats, draft readme and data dictionary
Linking Agent	Links related outputs, connect datasets/publications, PID assignment
Metadata Agent	Generates discipline aligned metadata, prepares metadata for deposit
Repository Agent	Discovers related datasets, recommend repositories (repository selection)
Provenance Agent	Records traceability, lineage
Validation Agent	Validates standards compliance, FAIR-alignment, governance constraints

9.5 Data Layer

The data layer manages all persistent storage, including dataset submissions, metadata records, policy configurations, audit logs and the knowledge sources used by the AI components (Table 8).

Table 8. Data layer components, with their descriptions and associated requirements, that manage persistent storage, metadata, lineage, and operational memory for the Data Director.

Component	Description/Requirement
Metadata Management	Data service capability that stores, governs, and exposes structured information about AI systems, data assets, models, agents, prompts, tools, and workflows to enable discoverability, governance, reuse, and traceability across the AI lifecycle
Lineage	Data service capability which records the end-to-end history, dependencies, and flow of data, prompts, models, and outputs across an AI system to provide full traceability from source to outcome
Storage	Foundational capability responsible for storing and retrieving all data, artefacts, and operational state required by AI models, agents, workflows, and supporting services
Immutable Storage	Storage capability designed to securely capture, preserve, and retain all critical AI system events, such as audit logs, user feedback, agent traces, and operational telemetry, in a tamper-evident, write-once, append-only format
Vector Database	Data service capability that stores, indexes, and retrieves high-dimensional vector representations of data to enable semantic search and Retrieval-Augmented Generation (RAG) in AI systems. Enables intelligent search and context injection by storing and retrieving meaning-based representations of enterprise knowledge to ground AI outputs in accurate, relevant information
Data Pipelines	Responsible for ingesting, conditioning, transforming, moving, and preparing data across systems while ensuring sensitive information can be securely masked, anonymised, or pseudonymised in accordance with governance, privacy, and compliance requirements.
NoSQL Document Store	NoSQL capability used as the operational memory layer for AI systems and agents. It stores semi-structured information such as conversation history (memory), user profiles, session state, workflow checkpoints, and derived “semantic memories” used for retrieval-augmented generation (RAG).

9.6 Platform and Integration Layer

The platform layer provides the computation, networking, identity and messaging infrastructure on which the Data Director runs. The integration layer manages connectivity with external systems (Table 9).

Table 9. Platform and integration components, with their descriptions and associated requirements, that manage computation, networking, identity and messaging infrastructure for the Data Director.

Component	Description/Requirement
API Management	Integration capability that provides a controlled, secure gateway for exposing, governing, and consuming APIs across internal and external services such as DataCite and Zenodo, etc.. It acts as the front door for all system-to-system communication, ensuring that APIs used by AI agents, workflows, applications, and external consumers are consistent, secure, and properly governed.
Identity Management	Foundational capability which is responsible for establishing, verifying, and governing the identities of users, applications, services, and AI agents, and controlling what they are allowed to access across the system. It acts as the trust and authentication backbone for authentication (confirming that an entity is who it claims to be) and authorisation (determines what an authenticated identity is permitted to do).
Security, Incident and Event Management	Capability which provides security observability and threat detection layer of an AI-enabled enterprise, providing real-time visibility into suspicious activity, policy violations, and potential security incidents across all system components, including AI agents, model runtimes, APIs, and data services. The service also supports real-time alerting and incident detection, triggering notifications when predefined rules or machine learning models detect anomalies, policy violations, or known attack signatures. These alerts can be routed to security teams, automated response systems, or incident workflows.
Operational Monitoring	Provides continuous visibility into the health, performance, availability, and reliability of applications, AI systems, workflows, and underlying infrastructure in real time.
Secret and Key Management	Responsible for securely storing, distributing, rotating, and controlling access to sensitive credentials such as API keys, encryption keys, certificates, passwords. It acts as the cryptographic trust and secret protection layer of an AI platform, ensuring that sensitive information is never exposed in code, configuration files, logs, or AI prompts, and is only accessed in a controlled and auditable way.

10. Evaluation Criteria

This section sets out how the Data Director Blueprint and any implementations of it will be evaluated. At v0.1, no reference implementation exists; this is a community specification at its first release. The evaluation framework is therefore defined here as a forward-looking commitment rather than an active assessment. It gives implementing organisations a clear target to design towards and gives the community a basis for assessing the Blueprint's impact as adoption grows.

10.1 Blueprint Quality Criteria

Quality criteria describe what a well-functioning implementation of the Data Director should demonstrate. The nine thematic areas below group the Blueprint's functional requirements (R), non-functional requirements (C), and architecture principles (P), reflecting the main dimensions across which an implementation can be assessed. Each area identifies what can be measured now and what cannot yet be measured.

Within each thematic area, criteria are distinguished by what the agent should do: its specific behaviours and outputs, and what the agent should be: its qualities and characteristics such as explainability, security, and accessibility.

*Within each area, items listed as **Cannot yet be measured** are tagged to indicate the nature of the gap. **[Agent/system development needed]** indicates that the gap requires further development of the agent, tooling, methodology, or community-agreed benchmarks. **[Human/user input needed]** indicates that the gap can only be closed through human judgement, structured user research, longitudinal adoption evidence, or direct involvement of community stakeholders*

1. Metadata and documentation quality | R2 · R3 · R4 · R5 · R7 · C15

A core output of the Data Director is high-quality metadata and supporting documentation. Two important caveats frame evaluation in this area. First, schema compliance alone is insufficient: values may be syntactically correct against a schema while being semantically wrong. Second, in domains where no controlled vocabulary, ontology, or prior FAIR practice exists, the tool's outputs will be generic, and the tool should explicitly flag this rather than presenting generic outputs as domain appropriate.

Can be measured now:

- **R2:** Ontology term coverage is measurable: the proportion of terms present in a draft metadata record against a specific vocabulary can be counted. Coverage should be interpreted alongside term appropriateness, as precise metadata may legitimately use only a small number of terms from a large vocabulary.
- **R3:** Vocabulary and ontology use verifiable against the FAIRsharing registry. For example, the MOMSI (Multi-Omics Metadata Standards Integration) FAIRsharing collection is confirmed useful as a cross-domain benchmarking reference.
- **R4:** If a given schema is adopted to describe data or metadata, the description must be valid against the chosen schema. Validation against common schemas including DataCite, Dublin Core, and discipline-specific standards should be testable through automated checking.
- **R7:** PID assignment workflow testable end-to-end against repository capabilities.

Cannot yet be measured:

- **[Agent/system development needed] C15:** No community-agreed benchmark exists for evaluating whether AI-generated metadata meets an acceptable curation standard.

- **[Human/user input needed] R5:** Whether data documentation is genuinely sufficient for data reuse by a third party requires human judgement and cannot be assessed automatically.
- **[Agent/system development needed]** Semantic correctness of metadata values: schema compliance does not confirm that values are meaningful in context, as demonstrated by the RDA Sample Type Working Group and RDA PIDINST Working Group, both of which identified difficulties aligning domain-specific concepts with DataCite.
- **[Agent/system development needed]** Standardised benchmarking datasets for cross-implementation comparison do not yet exist. Specifying such datasets, potentially covering Dublin Core plus a core set of key disciplinary metadata standards, would allow different Data Director implementations to be evaluated against each other over time.

2. FAIR and CARE principle alignment | P1 · R9 · R11 · C15

This area assesses whether the Data Director actively advances FAIR principles across all capabilities, appropriately applies CARE principles where relevant, retrofits metadata for already-published datasets, and links data to the broader research record. FAIR maturity can be assessed using existing community tools. No equivalent methodology currently exists for CARE. Where multiple FAIR evaluation tools produce differing scores, the tool should support users in selecting the most appropriate evaluation approach for their context rather than presenting all available scores simultaneously.

Can be measured now:

- **P1:** FAIR maturity scores assessable using community-recognised tools and comparable before and after implementation.
- **R11:** Dataset linking to publications, grants, software, and people via ORCID and ROR verifiable through automated testing.

Cannot yet be measured:

- **[Human/user input needed] P1 (CARE Principles):** No agreed evaluation methodology exists. Assessment requires direct involvement of Indigenous and community stakeholders and cannot be conducted through a checklist. CARE implementations are context-dependent and will require more than one profile.
- **[Agent/system development needed] C15:** No community-agreed benchmark exists for comparing AI-generated records against manually curated equivalents.

3. Interoperability and open standards | R6 · C5 · P2 · P6 · P7

This area assesses whether the Data Director operates using open, community-maintained standards, exposes interoperable APIs across instances, and defaults to open sharing. All criteria in this area are technically testable at deployment.

Can be measured now:

- **R6 and C5:** Metadata exchange across two or more Data Director instances without bespoke integration is testable; API conformance against a published specification is verifiable. The tool's ability to map metadata across multiple schemas, enabling deposit of the same dataset to different repositories without manual remapping by the researcher, is also testable.
- **P6:** All standards and dependencies in use are auditable against the open standards requirement.
- **P7:** A licence audit confirms the agent source code is publicly available under MIT or Apache 2.0.
- **P2:** A configuration audit confirms default sharing is open and any restrictions are accompanied by a documented legal basis.

4. Governance, provenance, and human oversight | R8 · R10 · C3 · C12 · C13 · P4 · P5

This area covers the traceability of agent actions, the enforceability of human oversight, and compliance with data governance and AI governance frameworks. Technical verification is possible for most criteria. Whether human oversight is genuinely meaningful in practice cannot be assessed technically and requires structured user research with researchers and data support professionals.

Can be measured now:

- **R10 and P5:** Provenance logging, audit trail completeness, and use of PROV-O verifiable through automated testing and independent audit.
- **R8:** DMP alignment verification at the point of data publication testable against machine-readable DMP outputs from compatible tools (DSW, Argos, DAMAP).
- **P4:** Technical presence of human approval controls and blocking of irreversible actions testable.
- **C12:** Enforceability of machine-readable governance policies verifiable through independent audit.

Cannot yet be measured:

- **[Human/user input needed] P4:** Whether human oversight is genuinely meaningful in practice requires structured user research.
- **[Agent/system development needed] C13:** Compliance coverage depends on which AI governance framework is adopted and varies by jurisdiction; full compliance cannot be assessed against a single universal standard.

5. Security, privacy, and sovereignty | C1 · C2 · C4 · P3 · P10

This area assesses whether security is built into the tool from the outset, whether privacy risks are detected and flagged rather than resolved autonomously, and whether data sovereignty requirements are respected across jurisdictions. All criteria are technically testable at deployment. Ongoing monitoring is additionally required beyond initial testing.

Can be measured now:

- **C1 and P3:** Penetration testing against a recognised framework such as NIST; encryption and role-based access control audit; audit log completeness verifiable.
- **C2:** Whether the tool correctly detects potential personal identifiers in data and metadata is testable. Whether the sharing workflow pauses and alerts the human reviewer upon detection is testable.
- **C4 and P10:** Data residency enforcement across configured jurisdictions testable; local and edge deployment options verifiable.

Cannot yet be measured:

- **[Agent/system development needed]** Ongoing monitoring is required to assess robustness against the tool being manipulated or redirected into taking actions outside its unintended scope. No single methodology currently covers this scenario comprehensively at the point of deployment.

6. Explainability | C14 · P8

This area assesses whether all AI-generated outputs are accompanied by human-readable explanations, whether explanations are genuinely understandable to non-technical users, and whether a log of agent thought processes is maintained for transparency and audit.

Can be measured now:

- **C14 and P8:** Presence of explanations alongside every AI output verifiable through automated testing.
- **C14:** Role governance of override rights testable against defined user roles.
- **P8:** Presence and completeness of thought process logs verifiable at deployment.

Cannot yet be measured:

- **C14:** Whether explanations are genuinely understandable to researchers, data support professionals, and non-technical users across different disciplines requires structured user research.

7. Accessibility and inclusive access | C6 · P12

This area assesses whether the Data Director meets accessibility standards across all interfaces and is genuinely designed to include rather than exclude researchers across communities, institution sizes, and levels of data expertise.

Can be measured now:

- **C6:** WCAG 2.1 compliance is independently auditable at deployment using recognised testing tools and assistive technology testing.

Cannot yet be measured:

- **[Human/user input needed] P12:** Whether the tool is genuinely inclusive across communities, institution sizes, and disciplines requires longitudinal evidence of adoption patterns and impact that does not exist at v0.1.

8. Operational performance | C7 · C8 · C9 · C10 · C11 · P11

This area assesses whether the Data Director meets reliability, availability, performance, and scalability requirements across deployment contexts. Most criteria are testable now. The 99.5% uptime target (C7) applies to the user-facing front end; internal components may have differential SLAs.

Can be measured now:

Uptime against the 99.5% target (C7), recovery from component failures without data loss (C8), 24/7 availability across time zones (C9), response time benchmarking and asynchronous processing (C10), load testing with growing volumes and concurrent users (C11), and deployment testing across edge, local, and cloud environments (P11) are all testable through standard infrastructure and performance testing at deployment.

These metrics apply to service-based deployments. Where the Data Director is deployed as a local or on-demand tool on personal or institutional hardware, infrastructure metrics such as uptime targets and concurrent user load testing may not be applicable. Implementers should assess which metrics are relevant to their deployment model and document any that are out of scope with justification.

Cannot yet be measured:

- **[Agent/system development needed] C7:** Implementing organisations must calibrate reliability targets to their deployment context. The 99.5% uptime target may not be realistic for locally hosted or resource-constrained deployments, such as those with limited computing power, IT support, or funding.

9. Impact, sustainability, and evaluation infrastructure | C16 · C17 · P9 · P13 · P14

This area covers two distinct dimensions that must not be conflated: the impact of the tool's outputs on researcher productivity and data quality, and the impact of the tool's own operation on the environment. Both must be monitored; neither yet has an agreed measurement methodology at v0.1. **Affordability (C17)** is identified as a future evaluation consideration in Section 10.3.

This area also covers the monitoring, logging, and feedback infrastructure needed to support ongoing evaluation of the tool's own performance and impact. Most mechanisms are technically verifiable at deployment.

Can be measured now:

- **C16:** Usage data collection and reporting verifiable at deployment. Monitoring dashboards, alert mechanisms, and feedback collection mechanisms all verifiable at deployment.
- **C17:** Configuration options for resource-constrained environments testable; AI processing cost architecture reviewable.
- **P9:** Local configuration points for schemas, repository preferences, and access controls testable.

- **P13:** Dependency viability assessable through review of licences, maintenance status, and community support. Testable by verifying that when dependencies are unavailable, unaffected workflow steps continue to function, the user is alerted, and work can be saved and resumed without data loss.

Cannot yet be measured:

- **[Human/user input needed] C16:** Whether usage data demonstrates genuine researcher impact requires longitudinal evidence. Automated quality scoring requires community-agreed thresholds that do not yet exist. User signal mechanisms require active researcher adoption before generating useful data.
- **[Agent/system development needed] C17:** No agreed methodology exists for assessing affordability across different institutional contexts and deployment models.
- **[Agent/system development needed] P14:** No agreed methodology exists for assessing environmental footprint across deployment models. Edge deployment is harder to monitor for environmental impact than centralised infrastructure, where active tracking tools can be deployed and usage and inference costs can be more reliably measured.

10.2 Evaluation Process

At v0.1, evaluation of any implementation is self-assessed by the implementing organisation against the quality criteria set out in section 10.1. Implementing organisations are encouraged to document which criteria they have addressed, record any limitations or gaps against the measurability ratings above, and share that evidence openly with the community in line with the Blueprint's Open First and Open-Source principles.

A formal evaluation process, including community review, independent verification, and a structured conformance framework, could be developed as a priority output of the v0.2 process, informed by feedback from early adopters during and after the 30-day community review period. Until that framework exists, the MUST designations in the functional specification (Section 7.1) define the minimum that any implementation claiming alignment with this Blueprint must satisfy.

10.3 Future Evaluation Considerations

Five areas of evaluation methodology remain underdeveloped in this version of the Blueprint. Each is identified for further development; several are carried forward explicitly in Section 13.2 (Roadmap for Future Iterations).

Community-agreed metadata quality benchmarks

The Blueprint requires AI-generated metadata to meet community curation standards (C15, R4), but no agreed threshold currently exists against which outputs can be assessed. Evaluation must extend beyond schema compliance to encompass semantic correctness: whether metadata values are accurate and meaningful in context. This is domain-specific and cannot be validated by machines alone. The Working Group has proposed a tiered quality

framework (gold, silver, bronze) as a direction for development in v0.2. Developing shared benchmarks will require structured collaboration across disciplinary communities, including through emerging frameworks such as [Assess-IF](#) and existing mechanisms such as FAIRsharing community collections and [FAIR Implementation Profiles \(FIPs\)](#).⁶⁰

CARE evaluation methodology

Existing community-recognised tools support FAIR maturity assessment. No equivalent approach for evaluating CARE principle application in practice currently exists, and assessment cannot be conducted through a checklist. It requires direct involvement of Indigenous and community stakeholders and must account for significant variation in CARE interpretation across communities and domains. As a minimum, the tool should identify when CARE principles are likely to apply and direct users to appropriate guidance; full assessment methodology remains a future development.

Longitudinal impact measurement framework

Whether the Data Director achieves equitable, sustained improvement in data publication practice cannot be determined at deployment (C16). A framework for measuring researcher uptake, equity of access, and the quality of data publication over time needs to be defined and agreed upon with the community. Standardised output formats, such as those produced by Assess-IF, can support longitudinal analyses once such a framework is in place.

Affordability and environmental impact assessment

The Blueprint requires the tool to be affordable across institutional contexts (C17) and to minimise its environmental footprint (P14), but specifies no measurement approach for either. No agreed methodology exists for assessing affordability across deployment contexts, and environmental monitoring is harder for edge deployments than for centralised infrastructure.

Formal conformance framework

Conformance at v0.1 is defined by satisfying the MUST requirements in the functional specification and is self-declared by implementing organisations. A community-governed conformance framework, including defined tiers, a peer review process, and an independent verification mechanism, could be developed as a priority output of the v0.2 process and piloted with early adopter implementations, replacing the current arrangement in which implementing organisations self-assess against the criteria without independent verification

11. Implementation Guidance

This section sets out practical starting points for organisations implementing the Data Director, drawn from reflections contributed by working group participants for each of the three implementer types identified in Section 2.2: research institutions, national research infrastructures, and commercial or non-profit vendors.

⁶⁰ <https://www.go-fair.org/how-to-go-fair/fair-implementation-profile/>

Some practicalities apply across all three implementer profiles, for example, considerations around cost transparency, interoperability, and conformance, while others are specific to a particular profile's role, resourcing, and relationship to the Blueprint. Where a consideration is genuinely cross-cutting, this is noted; otherwise, each subsection addresses the practicalities specific to that profile.

11.1 Cross-cutting Practicalities

The following practicalities were raised by more than one implementer profile and apply regardless of which profile an organisation belongs to.

Practical starting points:

- Find out what information about the tool's safety and security would be needed to satisfy a security team's assessment, including how data leakage is prevented.
- Establish how costs, including AI processing or token costs, would be estimated and how these would be passed on or absorbed (relevant whether estimating institutional running costs, infrastructure-wide sustainability costs, or end-user costs as a vendor).
- Understand how open standards and open-source commitments affect vendor lock-in and interoperability between implementations, instances, and institutions.
- Identify what counts as a valid implementation of the Data Director, including which capabilities are core to the specification versus optional or provided by separate agents, and how conformance with the Blueprint would be assessed or demonstrated.

11.2 Research Institutions

A research institution, in the context of this Blueprint, refers to an individual organisation, such as a university, research institute, or similar body, that conducts research and supports its researchers in producing and sharing research data. Research institutions are a primary deployment context for the Data Director, typically through their library, research support, or IT services, and are responsible for configuring the tool to reflect their own policies, systems, and governance arrangements.

The practical starting points below reflect considerations specific to deploying the Data Director within a single institution's existing systems, governance structures, and support arrangements.

Practical starting points:

- Identify which systems already deployed in the institution have known integration pathways with the Data Director and check against existing AI governance approval processes.
- Decide who will train research staff and whether institutional staff will provide day-to-day support; estimate how long it will take staff to get up to speed.
- Define a scoped minimum viable implementation, potentially as narrow as identifying relevant repositories and metadata standards or focusing first on open-access datasets before sensitive data.

- For smaller institutions, a first implementation focused on repository recommendation and basic metadata guidance for open datasets may be more achievable than a full deployment, with scope expanded once value and practical requirements are clearer.
- Plan how outputs feed back into institutional systems (e.g. CRIS, institutional repository) when a researcher deposits data in an external repository rather than the institution's own, so that a record is created locally pointing to where the data has been deposited.
- Identify what onboarding documentation will be needed for IT staff, administrators, and early users (beta-users) before deployment, for example, prerequisites, setup steps, and getting-started guidance, similar in style to vendor documentation such as learn.microsoft.com.⁶¹
- Clarify whether AI-assisted outputs can be claimed as research outputs or impact in national research assessment exercises (such as the UK's Research Excellence Framework (REF), Australia's Excellence in Research for Australia (ERA), or equivalent frameworks in other countries).
- For libraries, plan integration with university IT, Office of Research Development, and tools, such as research data repositories.

11.3 National Research Infrastructures

A national research infrastructure, in the context of this Blueprint, refers to a body operating at national (or potentially regional or international) scale that provides shared research data services, systems, or platforms used across multiple institutions within a country or research community. National research infrastructures differ from individual research institutions in that their decisions about adopting and configuring the Data Director have implications for the many institutions and researchers who rely on their services, and they must additionally consider cross-institutional and cross-jurisdictional factors.

The practical starting points below reflect considerations specific to deploying the Data Director at this larger scale, including business case development, integration with existing national services, and the legal and policy implications of acting on behalf of multiple institutions and users.

Practical starting points:

- Estimate what proportion of the national user base are candidate users, and identify existing infrastructure elements with known integration pathways, to build a business case.
- Confirm developer availability (in-house or elsewhere in the organisation) and establish a free or low-cost development and testing environment.
- Decide on the most important formats, vocabularies and schemas to support initially, rather than attempting broad coverage from the start.
- Plan the transition from the reference architecture to specific workflows, supported by documentation and training for researchers.

⁶¹ <https://learn.microsoft.com/en-gb/>

- Consider an SDK (Software Development Kit, a set of pre-built tools and code components that make it easier to build integrations with a system) or integration snippets as a minimum viable implementation.
- Be aware of the shift from 'user requests service' to 'Data Director requests service on behalf of user' and its implications for service design and user policy.

11.4 Commercial or Non-profit Vendors

A commercial or non-profit vendor, in the context of this Blueprint, refers to an organisation that develops, packages, hosts, or provides support for an implementation of the Data Director on behalf of research institutions or national research infrastructures, whether as a paid product or service (commercial) or as part of a not-for-profit or community-supported offering (non-profit).

The practical starting points below reflect considerations specific to building and offering an implementation of the Data Director to others, including business model and licensing questions, conformance, and avoiding duplication of existing implementations.

Practical starting points:

- Understand what business model could be built around the tool, whether it can be extended, and whether the licence constrains rights over your own IP.
- Identify what 'valid implementation' means and how conformance will be assessed for sign-off.

12. Blueprint Governance

The governance model defines how this Blueprint could be maintained, updated and kept aligned with community needs over time. Governance is the mechanism through which the community retains ownership of the specification and ensures it remains open, accurate and useful.

12.1 Ownership and Stewardship

This version of the Blueprint (v0.1) is created and owned by the RDA community, with the RDA positioned as an appropriate steward to keep the work open, transparent and neutral. Research institutions, infrastructures and vendors are positioned as contributors and implementers, rather than owners or governors. Ownership and governance of the Blueprint itself remain with the RDA community regardless of how the Blueprint is adopted. Separately, if the Blueprint develops into a standard, it would be useful to gather structured feedback from implementers who adopt it. This could include a forum where issues that become divisive among vendors could be discussed and reported back to the RDA community, helping to maintain the integrity and shareability of data for collaboration and metadata crosswalks.

12.2 Maintenance and Continuity

Stewardship could transition from the current time-limited Working Group to a standing maintenance function, or the Working Group could be re-formed on a time-limited basis for future iterations. Some form of ongoing maintenance presence within RDA structures may be needed, for example, a working group operating in a maintenance mode, or an interest group, with a rotating Chair and Co-Chair responsible for scheduling meetings, maintaining tools, working on improvements and updating policies, with rotation helping to address the risk of people losing interest over time. Either way, who carries out the work and how resources are provided would need to be addressed concretely, including how and where to recruit people to maintain the Blueprint within existing RDA structures.

Agentic models could support basic maintenance tasks, such as checking and updating links and flagging where referenced ontologies have changed. This refers to tooling that supports maintenance of the Blueprint document itself, such as checking external links or flagging changes to referenced ontologies and is separate from the Data Director as a deployed tool.

12.3 Platform and Transparency

The Blueprint could be maintained on [GitHub](#),⁶² allowing changes to be tracked, monitored and versioned, an approach supported by several contributors. A machine-readable, open format, such as Markdown, was suggested to reduce the risk of proprietary format changes affecting it over time, with GitHub potentially synced to [Zenodo](#)⁶³ or a similar platform for DOI minting and versioning. Wherever managed, the Blueprint should be openly accessible to all relevant stakeholders.

12.4 Review and Evolution

Feedback could be gathered from research institutions, infrastructures and vendors who have implemented the Blueprint, to inform updates over future review cycles, dependent on a future group or facilitated programme. There is potential to learn from how other organisations have approached the maintenance of similar technical solutions, including commercial editors such as [Springer](#),⁶⁴ [Elsevier](#)⁶⁵ and [MDPI](#)⁶⁶ who have implemented data repositories, AI developers with an interest in the scientific sector, and [W3C](#),⁶⁷ where examples exist of retro-compatibility being maintained between implementations and a specification. The Blueprint could evolve based on feedback from real-world implementations.

12.5 Scope and Sustainability

A question was raised as to whether this document is a seed point for implementers to build on as they see fit, or a technical specification that continues to be updated, the latter requiring considerably more ongoing work. Future iterations of the Blueprint would benefit from more

⁶² <https://github.com/ResearchDataAlliance>

⁶³ <https://zenodo.org/>

⁶⁴ <https://link.springer.com/>

⁶⁵ <https://www.elsevier.com/en-gb>

⁶⁶ <https://www.mdpi.com/>

⁶⁷ <https://www.w3.org/>

explicitly distinguishing between normative requirements that all implementations must satisfy, implementation guidance that supports but does not constrain implementation choices, and optional examples or future considerations that are illustrative or deferred. If the Blueprint continues as an ongoing activity, a sustainability model would be needed, including a concrete answer to how maintenance is resourced.

13. Open Questions and Future Iterations

This section documents the current boundaries of the Blueprint. Section 13.1 sets out design questions that remain unresolved in this version, for consideration as the specification develops and section 13.2 identifies considerations deferred to future iterations.

13.1 Unresolved Design Questions

Table 10 presents design questions identified during the Blueprint programme that remain unresolved in this version, for consideration as the specification develops.

Table 10. Open questions raised during the development of this Blueprint that remain unresolved in this iteration and require further community discussion, governance decisions, or design work before they can be closed.

Theme	Question	Notes
FAIR evaluation	Should the Data Director support evaluation of already-published digital objects (e.g. via FAIRassist), in addition to metadata produced during the publication workflow?	The community has indicated this would be valuable. Open questions remain about how the tool connects to external evaluation services post-publication, how results are surfaced, and how findings are recorded in the provenance chain.
FAIR evaluation	Should FAIR be assessed at the FIP level rather than FDO level, and separately at repository, dataset, and file levels, given the current architecture conflates these?	Structural issue affecting the uses table (Section 4.1) and reference architecture (Section 9).
Openness and sovereignty	Should P2 (Open First) be reframed as “as open as possible, as sovereign, protected, and accountable as necessary,” conditioned on human rights, consent, and community authority?	Raised as critical for international projects and genuine inclusivity, to avoid widening the digital divide. Needs resolving before national-level deployment.
Openness and sovereignty	Should computation be brought to the data, rather than data moved to computation, where residency requirements prevent central processing?	A recognised pattern in sensitive data environments, but the practical means of achieving it depends on institutional infrastructure and is not specified. Implementers under strict sovereignty or residency constraints should treat this as requiring institutional design.

AI safety, resilience and evaluation governance	<p>How should the safety, transparency and resilience of the Data Director be assessed, reported and governed over time, including where humans remain in the loop, what outputs must always be reviewed, and how escalation and suspension should work?</p>	<p>C13 and C14 set requirements for responsible AI behaviour, but how the tool is audited against emerging governance frameworks, monitored in production, and trusted as capabilities evolve remains open. Identified as the key community governance question; decisions are expected to need review as agent competence grows, with the appropriate level of oversight likely determined by the funder or institution based on resourcing and workload.</p>
Interoperability	<p>How should system-level research data exchange and reuse across institutional boundaries be supported, beyond metadata format alignment (R6)?</p>	<p>Encompasses authentication, authorisation, data transfer protocols, and institutional data governance agreements. Designing a system-level framework is beyond this Blueprint's scope; a question for future iterations and related infrastructure initiatives.</p>
Agent memory and model improvement	<p>Should researcher query data be used to improve the underlying AI model over time?</p>	<p>Involves a balance between privacy and quality improvement; the answer will vary by institution, jurisdiction, and governance framework. Whatever decision is made must be documented, communicated to users, and reviewed periodically.</p>
Scope and data types	<p>Is the Data Director workflow adaptable to dynamic data; continuously produced, high-volume outputs such as environmental sensor streams or geophysical feeds, or does it require adjustment for research infrastructure contexts where data is not static?</p>	<p>The Blueprint is designed primarily with static datasets in mind. Whether the workflow applies as-is or requires adaptation for dynamic data production environments is an open question for future community consideration.</p>

13.2 Roadmap for Future Iterations

This section sets out items deferred from v0.1, representing areas of strong community interest that are either technically premature, outside the current scope, or dependent on decisions not yet made.

Evaluation and quality

- Permit a first-pass review by another AI agent before human review, to manage scale, with the ability for a human reviewer to revert to a prior agent version if unsatisfied.
- Include measurement of hallucination levels in AI-generated metadata as part of quality evaluation under C15, dependent on the state of the art for hallucination detection.

Scope and lifecycle

- Engage researchers at the point of data collection and project planning, not only at publication; automating engagement at the point of award was suggested as a way to address this.

Workflow and architecture

- Add pipeline extensibility as a named architectural principle, enabling optional stages such as anonymisation, format conversion, and compliance checking.
- Make access procedures interoperable between Data Director instances (per C5), with technical mechanisms to harmonise across a federation of agents, while recognising that access policies (authorisation rules) cannot be identical across organisations.
- Extend R12 (Discover related datasets) to include detection of potential duplicate datasets, particularly relevant in large consortia where the same dataset may be deposited more than once by different teams.
- Consider whether the linking of datasets to related research outputs via R11 could be extended to surface technology transfer pathways; routes through which research outputs may have commercial, industrial or societal application beyond the original research context, noting that identifying such pathways should remain a conclusion the researcher draws from seeing related outputs, rather than an explicit suggestion made by the tool.

Sovereignty and data residency

- Clarify how data sovereignty should be enforced in practice when AI processing, logs, temporary files, and backups may involve external services or cross-border infrastructure, with national research agencies involved in RDA potentially solicited for input.

Accessibility and language

- Support multilingual metadata records (C6) and define how the tool should handle datasets containing multiple languages. One view is to leave this as a configurable, agentic capability rather than a default requirement, important for inclusivity. Open questions include the effect on performance, minimum language coverage, and whether outputs should default to an English-language version alongside others.

Naming and terminology

- Review whether the name 'Data Director' risks implying the tool directs, authorises, or takes responsibility for decisions. Possible alternatives considered include Research Data Assistant, Data Concierge, Data Curator, and Data Curator Agentic AI, with a suggestion that adding 'agent' to the name could help convey its role without further explanation.
- Consider replacing the term 'publication' throughout with a less ambiguous term such as 'research output' or 'digital research object', affecting precision throughout the document.

14. Appendices

Appendix A: Glossary of Terms

Terms are listed alphabetically. This glossary covers terminology, acronyms, and external standards used in the Blueprint without an in-text explanation. Sources are provided for each entry and hyperlinked where available. Definitions have been prepared with AI assistance and represent a combination of direct quotations, paraphrases, and working definitions informed by the cited sources. Where precise or authoritative definitions are required, readers are encouraged to consult the cited source directly.

Term	Definition (with source)
Anonymisation and Pseudonymisation	Anonymisation removes or alters personal data so that individuals can no longer be identified, even indirectly. Pseudonymisation replaces identifying details with artificial identifiers, so that re-identification remains possible only with additional information held separately and securely. (Source: Information Commissioner's Office)
API (Application Programming Interface)	A defined set of rules and protocols that allows different software systems to communicate with one another, for example, enabling the Data Director to exchange information with a repository, registry or identifier service. (Source: MDN Web Docs)
CARE Principles	Collective Benefit, Authority to Control, Responsibility and Ethics: principles developed to ensure Indigenous Peoples retain control over how data about their communities, lands and knowledge is collected, used and shared, complementing the FAIR principles. (Source: Carroll, S.R. et al. (2020). The CARE Principles for Indigenous Data Governance. Data Science Journal, 19(1), p.43. https://doi.org/10.5334/dsj-2020-043)
CRedit (Contributor Roles Taxonomy)	A standardised taxonomy of contributor roles used to describe the specific nature of each person's contribution to a research output. Referenced as the basis for recording human roles in provenance records. (Source: CRedit)
CRIS (Current Research Information System)	An institutional database system used to capture, manage and report on information about research activity, outputs, funding and personnel. (Source: euroCRIS)
Crosswalk	A mapping that shows how fields, terms or elements in one metadata schema correspond to equivalent fields in another, enabling records created for one system to be understood and reused in another. (Source: Library of Congress)
Data Dictionary	A document that lists and describes the variables, fields and codes used in a dataset, including their definitions, units and permitted values, helping others understand and reuse the data. (Source: The Turing Way)
Deposit (Research Data)	The act of submitting a research dataset to a repository for long-term storage, preservation, and access. Distinct from publication in that deposit is the researcher's act of submitting data, while making data publicly available is typically carried out by the repository following the deposit. The Data Director supports researchers in preparing for and completing deposits; it does not itself publish or make data publicly available. (Source: OECD Principles and Guidelines for Access to Research Data from Public Funding , 2007; see also COAR Controlled Vocabulary for Resource Type Genres)
Digital Object Identifier (DOI)	A persistent, internationally standardised identifier used to uniquely and permanently identify a digital object such as a journal article, dataset or report, ensuring it remains findable even if its web location changes. (Source: DOI Foundation ; see also ISO 26324:2022, Information and documentation - Digital object identifier system)"

Dublin Core	A widely used, general-purpose metadata standard consisting of a small set of core elements, such as title, creator, date, and subject, for describing digital and physical resources, designed to support resource discovery across domains. (Source: Dublin Core Metadata Initiative)
Edge Deployment	An approach to deploying computing infrastructure close to where data is generated or used, such as on a local server or device, rather than relying on centralised cloud data centres, reducing latency and supporting data residency requirements. (Source: ISO/IEC TR 23188:2020, Information technology - Cloud computing - Edge computing landscape , ISO/IEC Joint Technical Committee 1, Subcommittee 38)
Embargo	The set period of time by which investigators retain exclusive rights to their dataset while it has been deposited in a repository but before it is made publicly available. (Source: SPARC Portal documentation, https://docs.sparc.science/docs/embargoed-data , adapted for general use)
FAIR Digital Object (FDO) framework	A machine-actionable digital object that bundles data, metadata, type information and a persistent identifier together, enabling software systems to automatically interpret, validate and act on its content and structure. (Source: FAIR Digital Object Forum)
FAIR Principles	Findable, Accessible, Interoperable and Reusable: a set of guiding principles for making research data, software and other digital objects usable by both humans and machines, and a cornerstone of open science policy. (Source: Wilkinson, M.D. et al. (2016). The FAIR Guiding Principles for scientific data management and stewardship. Scientific Data, 3, 160018. https://doi.org/10.1038/sdata.2016.18)
Federation	A model for connecting multiple independent research data systems, repositories or infrastructure nodes so that resources, services and datasets can be discovered, accessed and used across them without being centralised or merged into a single system. Each participant retains ownership and control of their own resources. (Source: European Open Science Cloud (EOSC) Architecture Documentation . 'The European Open Science Cloud (EOSC) is a federated system designed to enable seamless sharing, discovery, and access to research resources across Europe. It consists of interconnected EOSC Nodes, which can be national, regional, or thematic infrastructures that provide services, datasets, and computational resources to the research community.'
GUPRI (Globally Unique, Persistent and Resolvable Identifier)	An identifier that is globally unique, remains valid and resolvable over time, and can be used by humans and machines to access information about the digital object it identifies. As described in the EOSC PID Policy, a resource that is FAIR must be assigned an identifier that meets all three criteria. (Source: EOSC FAIR Working Group and Architecture Working Group. A Persistent Identifier (PID) policy for the European Open Science Cloud . European Commission Publications Office, October 2020. See also FAIRsharing documentation on GUPRIs .)
Hallucination (AI)	A term used to describe instances where an AI system generates content that is confidently stated but factually incorrect, erroneous, or not grounded in verifiable evidence. The term is contested: the output may be consistent with patterns in training data while still being factually wrong or insufficiently evidenced for the purposes for which it is used. NIST AI 600-1 uses the more precise term confabulation to describe the production of confidently stated but erroneous or false content. (Source: NIST AI 600-1, Artificial Intelligence Risk Management Framework: Generative Artificial Intelligence Profile , National Institute of Standards and Technology, 2024)
Handle System	A system for assigning and resolving persistent digital identifiers, known as Handles, which provide a stable reference to a digital resource regardless of changes to its location. (Source: Handle.net Registry)

Infrastructure as Code (IaC)	The practice of defining and managing computing infrastructure, such as servers, networks and storage, through version-controlled configuration files rather than manual processes, enabling consistent, repeatable and auditable deployment. (Source: Microsoft Learn)
IRB / REC (Institutional Review Board / Research Ethics Committee)	A formally constituted committee responsible for reviewing and approving the ethical aspects of research involving human participants, animals or sensitive data, ensuring that research meets internationally recognised ethical standards before it proceeds. Known as an Institutional Review Board (IRB) in the United States and a Research Ethics Committee (REC) in many other countries. (Source: World Health Organization. Standards and operational guidance for ethics review of health-related research with human participants . WHO, 2011. See also CIOMS International Ethical Guidelines for Health-related Research Involving Humans , 2016, developed in collaboration with WHO)
JSON Schema	A vocabulary for describing the structure, content and validation rules of JSON (JavaScript Object Notation) data, used to check that a metadata record or dataset conforms to an expected format. (Source: JSON Schema)
Knowledge Graph	A structured network of entities, such as datasets, people, organisations or concepts, and the relationships between them, represented in a machine-readable form that supports search, discovery and reasoning. (Source: W3C)
Machine-readable	Describes information that is structured and formatted so that it can be automatically read, parsed and processed by computer systems, as distinct from formats intended primarily for human reading. (Source: UNESCO Recommendation on Open Science)
Metadata	Structured information that describes, explains, or locates a resource, enabling it to be found, understood, and reused. (Source: Dublin Core Metadata Initiative / ISO 15836). For the purposes of this Blueprint, metadata is understood in its full scope, encompassing not only discovery fields but variable definitions, data structure, provenance, quality assessments, access conditions, and applicable governance constraints. See Section 3.1.1.
MLCommons Croissant	A community-developed, open metadata format for describing machine learning datasets in a structured, machine-readable way, designed to make datasets easier to find, load and use across different tools. (Source: MLCommons)
Model Card / Dataset Card	A short, structured document accompanying a machine learning model or dataset that records its intended use, composition, known limitations and evaluation results, intended to improve transparency and reusability. (Source: Hugging Face)
Nanopublication	A small, citable unit of scientific knowledge expressed as a single machine-readable assertion together with information about its provenance and supporting evidence, designed to be linked together to form larger knowledge graphs. (Source: Nanopublications)
NIST (National Institute of Standards and Technology)	A US government agency that develops widely adopted standards and frameworks, including in cybersecurity and AI risk management, often used as a reference baseline for security practice internationally. (Source: NIST)
ODRL (Open Digital Rights Language)	A standardised, machine-readable language for expressing permissions, prohibitions and obligations relating to the use of digital content and data, such as licensing or access conditions. (Source: W3C)
Ontology	A formal, machine-readable representation of the concepts within a domain and the relationships between them, used to support consistent classification, search and integration of data across systems. (Source: W3C Semantic Web)

Open Science	A broad movement and set of practices aimed at making the processes, outputs and tools of research, including publications, data, methods and software, openly available, accessible and reusable for the benefit of science and society. (Source: UNESCO Recommendation on Open Science)
Open-Source Licence (e.g. MIT, Apache 2.0)	A type of software licence that grants users the right to freely use, modify and distribute software. The MIT licence is short and highly permissive; the Apache 2.0 licence is similarly permissive but additionally includes an express patent licence. (Source: Open Source Initiative)
ORCID	A free, persistent identifier that researchers use throughout their career to distinguish themselves from others and to link their identity to their research outputs, affiliations and activities. (Source: ORCID)
Persistent Identifier (PID)	An identifier assigned to a digital resource, person or organisation that is globally unique, remains valid over time, and can be resolved by both humans and machines to locate or access the identified entity. A resource that is FAIR must be assigned a PID that meets these three criteria of global uniqueness, persistence, and resolvability. (Source: EOSC FAIR Working Group and Architecture Working Group. A Persistent Identifier (PID) policy for the European Open Science Cloud . European Commission Publications Office, October 2020)
PROV-O (PROV Ontology)	A standard vocabulary for recording provenance information in a machine-readable form, describing the entities, activities and people involved in producing, influencing or delivering a piece of data. (Source: W3C)
RAG (Retrieval-Augmented Generation)	An AI technique in which a language model's response is grounded by first retrieving relevant information from an external knowledge source, such as a document collection or database, and supplying it as additional context before generating an answer. (Source: AWS — What is RAG?)
Repository (Research Data Repository)	An online platform for storing, preserving, describing and providing access to research data and other outputs over the long term, often assigning persistent identifiers and supporting discovery through metadata. (Source: The Turing Way)
ROR (Research Organization Registry)	An open, community-led registry providing persistent identifiers for research-performing organisations, used to unambiguously link research outputs, funders and institutions. (Source: ROR)
SHACL (Shapes Constraint Language)	A standard language for describing and validating the structure of data, used to check that metadata conforms to a defined shape or set of rules. (Source: W3C)
WCAG (Web Content Accessibility Guidelines)	An internationally recognised set of guidelines defining how to make web content more accessible to people with disabilities. (Source: W3C Web Accessibility Initiative)
XSD (XML Schema Definition)	A standard for describing the structure, content and permitted data types of an XML document, used to validate that an XML file conforms to an expected format. (Source: W3C)

Appendix B: Related Resources

The resources below are drawn from across the Blueprint and related working group discussions. They are grouped by theme, with the resource title hyperlinked to its primary URL. Entries with no confirmed URL in the source material are shown in bold without a link, with this noted in the description.

RDA, Governance and Consultation

Resource	Description and source section
RDA About	General information about the RDA.
RDA Guiding Principles	The Guiding Principles with which the Blueprint's openness and governance are designed to comply.
Data Director Agentic AI Blueprint Working Group	Working Group activity page for this Blueprint.
Global Community Priorities for Agentic AI Development	Outputs of the RDA-Microsoft global community consultation that selected the Data Director as the Blueprint's focus.
How RDA Works	Describes the RDA's mission, vision and guiding principles, against which the Data Director was assessed for alignment.
RDA Recommendations and Outputs process	The formal RDA process under which this Blueprint is produced.
RDA Private Sector Engagement Framework v1.0	Framework under which Microsoft's facilitation of the Working Group was provided.
RDA Guidance on AI Tools Usage	Guidance followed for the use of Claude in drafting and analysis throughout the Blueprint.
RDA Decentralised Action Repository recommendations	RDA recommendations that could inspire a decentralised action repository supporting full agent traceability.
RDA maDMP Working Group	Maintains a fuller list of DMP tools compatible with machine-actionable Data Management Plans.

FAIR, CARE and Indigenous Data

Resource	Description and source section
FAIR Principles	The Findable, Accessible, Interoperable, Reusable principles underpinning the Data Director.
Comprehensive FAIR Principles paper (Wilkinson et al.)	Foundational publication defining the FAIR principles.
CARE Principles for Indigenous Data Governance	Collective Benefit, Authority to Control, Responsibility and Ethics: principles for Indigenous data governance applied alongside FAIR.

Global Indigenous Data Alliance (GIDA)	Community standards body providing explicit CARE guidance for Indigenous data.
Local Contexts	External system identified as relevant for meaningful engagement with CARE principles.
GO FAIR	Community resource referenced for domain-appropriate standards recommendations and the FAIR principles.

Persistent Identifiers and Identifier Frameworks

Resource	Description and source section
DataCite	Primary membership model for institutional DOI minting for deposited datasets.
FAIRsharing identifier schema catalogue	Catalogues identifier schemas by persistence, resolvability and global uniqueness per the EOSC PID definition.
EOSC PID definition	Community-agreed definition of persistent identifier properties referenced in the Blueprint's PID requirements.
W3C Decentralised Identifiers (DID) v1.1	Standard for decentralised identifiers, proposed as a basis for both dataset and agent identity.
ORCID	Persistent identifier for researchers; identified as a linking target for related research outputs.
ROR (Research Organisation Registry)	Persistent identifier for research organisations; identified as a linking target for related research outputs.

Repository and Standards Registries

Resource	Description and source section
FAIRsharing	Curated registry of community-defined standards, databases, policies, identifier schemas and reporting guidelines across disciplines.
FAIRsharing MCP server	MCP server enabling direct integration of FAIRsharing with the Data Director.
re3data	Registry of research data repositories with filtering by subject, access type and deposit conditions.

Cornell Data Storage Finder	Open-source tool for localised repository and storage recommendations; suggested as a possible reference implementation.
CoreTrustSeal	Repository certification and trustworthiness scheme; free certification options should be prioritised over paid alternatives.

FAIR Evaluation and Quality Tools

Resource	Description and source section
FAIRassist.org	Full list of community FAIR evaluation tools, including those with OSTRails framework integration.
F-UJI	Automated FAIR data assessment tool.
FAIR EVA	FAIR Evaluation and Validation Assistant.
FAIR Checker	FAIR data assessment tool.
FAIR Champion	FAIR data assessment tool.
FAIRshake	Metadata quality evaluation framework for assessing AI-generated record quality against curation standards; relevant to C15.
CESSDA FAIR benchmark	Example of a community-defined FAIR benchmark registered in FAIRsharing.
MOMSI FAIRsharing collection	Cross-domain benchmarking reference covering domain-specific and universal standards; confirmed useful for evaluating R3 and R4.
Assess-IF framework (OSTRails)	Framework for community-specific FAIR benchmark definition and longitudinal analysis; implements the open framework for FAIR evaluation scores (w3id.org/fttr#).
FIPs (FAIR Implementation Profiles)	Community-defined FAIR benchmarks registrable in FAIRsharing.

Metadata Standards, Schemas and Validation Tools

Resource	Description and source section
RDA kernel metadata recommendations	Defines essential domain-agnostic metadata fields.

FAIR2 metadata specification	Community-driven metadata specification built on MLCommons Croissant.
Ontology Lookup Service (OLS)	Ontology repository providing controlled vocabularies across disciplines.
NFDI Terminology Service	Terminology service providing controlled vocabularies for research communities.
OntoPortal	Domain-specific ontology repository platform.
BioPortal	Domain-specific ontology repository (biomedical).
AgroPortal	Domain-specific ontology repository (agronomy).
RO-Crate validator	Schema validation tool for metadata quality checks.
Frictionless Data	Schema validation tool for metadata quality checks.
JSON Schema	Schema validation standard applicable to metadata validation.
XSD (XML Schema Definition)	XML schema validation reference.
SHACL (Shapes Constraint Language)	Schema validation tool for metadata quality checks.
Croissant Validator / Croissant ML	Mature, generic metadata format and validator with an active community; proposed as a pilot format, supported by Croissant Baker and Croissant Miner.
CEDAR Workbench	Metadata creation and validation tool.
Skosmos	Vocabulary publishing and validation tool.
FOOPS!	Ontology pitfall scanner for vocabulary validation.
PROV-O	Candidate standard for provenance recording with archiving option, identified for tracking AI agent actions and Full Traceability.
eScience 2025 provenance paper	Community reference for PROV-O implementation in agentic AI research workflows.
ODRL (Open Digital Rights Language)	Policy expression language referenced for machine-readable governance policies; proposed as a FAIR object with identifiers. Demo: odrl.dev.codata.org/demo .

FAIR Digital Object (FDO) framework	Supports interoperability through typed, validated objects, enabling metadata at domain-agnostic, domain-specific and application-specific levels simultaneously.
Cross-Domain Interoperability Framework (CDIF)	Supports cross-domain metadata mapping; relevant to interoperability requirements.

Data Management Plan (DMP) Tools

Resource	Description and source section
Data Stewardship Wizard (DSW)	Exposes DMP commitments in a computationally amenable form; confirmed compatible with the Data Director.
Argos	DMP tool with machine-readable output; confirmed compatible with the Data Director.
DAMAP	DMP tool with machine-readable output; confirmed compatible with the Data Director.

AI Governance, Safety and Security

Resource	Description and source section
Australian Cyber Security Centre - careful adoption of agentic AI services	Guidance on careful adoption of agentic AI services, directly relevant to the agentic AI security baseline.
Agentic Automation Canvas	Structured tool for honest planning of agentic platforms covering benefits, risks, human oversight and evaluation. Companion paper: arxiv.org/abs/2602.15090 .
Agentic AI expectations vs realised value (companion paper)	Companion paper on the gap between agentic AI expectations and realised value.
Policy Cards	Deployment-level governance tool that feeds from the Agentic Automation Canvas.
ACHPR Draft African Guidelines on Data Access and Human Rights	Draft guidelines on promoting and harnessing data access for human rights and sustainable development; no URL given in source material. Relevant to legal/sovereignty prerequisites and Open First.

<p>EOSC-Future / RDA ENVISAGE & PILOT Principles for AI</p>	<p>Principles under development relevant to a potential future Ethics by Design architecture principle.</p>
---	---

Environmental Impact

Resource	Description and source section
<p>OSAI Guidelines on AI environmental impact</p>	<p>Includes Frugal AI principles and Green AI techniques relevant to Low Impact (P14).</p>
<p>OSAI AI ecosystem landscaping</p>	<p>Calculators and tooling for environmental monitoring relevant to Low Impact (P14).</p>
<p>OECD AI Footprint</p>	<p>Environmental footprint resource relevant to Low Impact (P14).</p>
<p>FAS energy footprint measurement</p>	<p>Energy footprint measurement resource relevant to Low Impact (P14).</p>
<p>Nature - AI environmental impact</p>	<p>Citable reference for environmental impact considerations relevant to Low Impact (P14).</p>

Linking, Discovery and Impact Measurement

Resource (linked)	Description and source section
<p>COMET initiative</p>	<p>Supports community-curated enrichment of PID metadata for existing outputs, directly aligned with the retrospective scope of metadata remediation.</p>
<p>DataCite MCP</p>	<p>MCP server for dataset discovery, in active development at time of writing; relevant to surfacing related datasets.</p>
<p>MicroPublication Biology</p>	<p>Referenced as a potentially relevant resource for linking data to minimal research outputs.</p>
<p>CSIRO Impact Measurement Framework</p>	<p>Includes lead indicators for slow-materialising impacts; proposed as a model for longitudinal evaluation of the Data Director's impact.</p>

Comparator Models for Maintenance and Review

Resource	Description and source section
GitHub (RDA organisation)	Proposed platform for maintaining the Blueprint with tracked, monitored and versioned changes.
Zenodo	Proposed platform, potentially synced with GitHub, for DOI minting and versioning of the Blueprint.
Springer	Commercial publisher with experience implementing data repositories, referenced as a maintenance comparator.
Elsevier	Commercial publisher with experience implementing data repositories, referenced as a maintenance comparator.
MDPI	Commercial publisher with experience implementing data repositories, referenced as a maintenance comparator.
W3C	Referenced as an example where retro-compatibility has been maintained between implementations and a specification.

Working Group Tools and Documentation

Resource (linked)	Description and source section
Zoom	Used to facilitate Working Group sessions across global time zones.
Google Docs	Used for collaborative session notes.
Draw.io	Used for collaborative diagramming.
Miro	Used for collaborative working and generating process flow and architecture diagrams.
Claude.ai	Used to generate final process flow images, with Claude Sonnet 4.5 (Pro) used for data analysis and writing assistance throughout.
learn.microsoft.com	Referenced as a documentation style example for institutional onboarding materials.

Appendix C: Acknowledgements

First and foremost, thank you to all the active members of the working group for their dedication and commitment to the programme over the eight-week period (April - June 2026). Thank you for attending the many sessions held at different times across the globe, spanning different regions and time zones, and for your openness, willingness, and enthusiasm to share knowledge each week. It is this collective effort that has led to the generation of such breadth and depth of content, culminating in this comprehensive draft Blueprint.

Special thanks to Connie Clare, RDA Community Development Manager and facilitator of the working group for her passion, drive and dedication to this activity. Much gratitude to Hilary Hanahoe (RDA Secretary General) and Marcy Collinson (Director, Worldwide Academic Research at Microsoft) for bringing this public-private partnership between RDA and Microsoft to fruition, and to the software engineers Benjamin Wright-Jones and George Earl for always being on hand to support the planning and execution of the sessions. Their time, expertise, and effort, particularly in developing the reference architecture in Section 9 has been invaluable. Grateful thanks also extend to Trish Radotic (RDA Secretariat) for serving as a sounding board and providing logistical support throughout the process.

Blueprint Authors

As of June 2026, the [Data Director Agentic AI Blueprint Working Group](#) comprises 94 members. A significant number of these individuals have actively contributed to the development of this Blueprint by attending sessions, participating in discussions, and collaboratively developing the Blueprint through shared notes, Miro boards, and direct editorial contributions.

No.	First Name	Last Name	Role	Affiliation	Country	ORCID
1	Maricela	Abarca	Data Curator	Stanford University Libraries	United States	0000-0002-0890-8887
2	Lindsey	Anderson	Computational Scientist	Pacific Northwest National Laboratory (PNNL)	United States	0000-0002-8741-7823
3	Julieta	Arancio	Senior Advisor	ORCA	Germany	
4	Amir	Aryani	Head of AI Group	Swinburne University of Technology	Australia	
5	Rossella	Aversa	Head of Department	Karlsruhe Institute of Technology (KIT)	Germany	0000-0003-2534-0063
6	Amira	Azizan	Data Steward	Centre National de la Recherche Scientifique (CNRS) / Data Terra	France	
7	Sudhanshu	Bajpai	Librarian	Indian Institute of Technology BHU (IIT BHU)	India	
8	Laure	Berti-Equille	Research Director	Institut de Recherche pour le Développement (IRD)	France	
9	Paty	Buendia	Technical Director	Lifetime Omics	France / United States	
10	James	Cannon	Research HCP and Storage Engineer	University of Michigan	United States	

11	Connie	Clare	Community Development Manager	Research Data Alliance (RDA)	United Kingdom	0000-0002-4369-196X
12	Francis P.	Crawley	Executive Director	EOSC-Future/RDA Artificial Intelligence & Data Visitation (AIDV) WG; HGP2; GCPA; SIDCER	Belgium	0000-0002-6893-5916
13	Mohamed	Drira	Associate Professor	Saint Mary's University	Canada	0000-0002-8952-8148
14	George	Earl	Software Engineer	Microsoft	Canada	
15	Joe	Edgerton	Research Data Management Librarian	University of Virginia	United States	0009-0000-1292-6642
16	Lars	Eklund	Domain Specialist Sensitive data	National Bioinformatics Infrastructure Sweden (NBIS) / Uppsala Multidisciplinary Center for Advanced Computational Science (UPPMAX) / Swedish National Data Service (SND), Uppsala University	Sweden	
17	Eva Eleanora	Ferradosa	Researcher	Delft University of Technology	Netherlands	
18	Roberta	Ferretti	Researcher	Italian National Research Council – Institute of Marine Engineering (CNR-INM)	Italy	0000-0002-1985-2145
19	Angela	Fuentes Pardo	Data Steward	Science for Life Laboratory (SciLifeLab) Data Centre	Sweden	
20	Marcelo	Garcia	Library Systems & Integrated Services Lead	King Abdullah University of Science and Technology (KAUST)	Saudi Arabia	
21	Mohit	Garg	Assistant Librarian	Indian Institute of Technology Delhi (IIT Delhi)	India	
22	Su Nee	Goh	Librarian	Nanyang Technological University	Singapore	
23	Cristina	Gonzalez	AI Data Specialist	Senscience	Switzerland	
24	Ritu	Gupta	Business Process Consultant	University of Michigan	United States	
25	Joe	Heffer	Senior Research Data Engineer	University of Sheffield	United Kingdom	0000-0001-8733-1117
26	Simon	Hodson	Executive Director	Committee on Data of the International Science Council (CODATA)	France	
27	Jenn	Huck	Associate Director, Research Data Services	University of Virginia	United States	0000-0003-4945-839X
28	Maximilian	Inckmann	Researcher	Karlsruhe Institute of Technology (KIT)	Germany	0009-0005-2800-4833
29	Agate	Jarmakoviča	Data Steward	Rīgas Stradiņš University (RSU)	Latvia	
30	Seonyoung	Kim	Senior Support Scientist	Washington University	United States	
31	Jens	Klump	Group Leader	Commonwealth Scientific and Industrial Research Organisation (CSIRO) / AuScope	Australia	

32	Piotr	Krajewski	Data Steward	Gdańsk University of Technology (Gdańsk Tech)	Poland	
33	Allyson	Lister	FAIRsharing Content & Community Lead	University of Oxford	United Kingdom	0000-0002-7702-4495
34	Sebastian	Lobentanzer	Principal Investigator	Helmholtz Munich	Germany	
35	Carolina	Mazza	Researcher	Universidad Nacional de La Plata (UNLP)	Argentina	
36	Marina	McGale	Technical Manager	Australian Data Archive (ADA)	Australia	
37	Ethel	Mendocilla Sato	Data Steward/Student	University of Lausanne (UNIL) / École Polytechnique Fédérale de Lausanne (EPFL)	Switzerland	0000-0003-0339-3535
38	Natalie	Meyers	AI Researcher in Residence/Research Specialist	Association of Research Libraries (ARL) / Coalition for Networked Information (CNI) / San Diego Supercomputer Center (SDSC)	United States	
39	Natalia	Mikołajczak	Open Science Specialist, Open Science Coordinator at WUT, Data Steward	Warsaw University of Technology, Main Library	Poland	0000-0002-6605-5368
40	Naeem	Muhammad	Research Data Manager	KU Leuven	Belgium	
41	Hari	Nepal	Researcher	University of Porto	Portugal	0000-0003-3078-3846
42	Anh	Nguyet Vu	Bioinformatics Engineer	Sage Bionetworks	United States	
43	Amy	Nurnberger	Head, Data Management Services	Massachusetts Institute of Technology (MIT)	United States	
44	Wolmar	Nyberg Åkerström	Bioinformatician / Data Steward	European Life Sciences Infrastructure for Biological Information (ELIXIR) Sweden / National Bioinformatics Infrastructure Sweden (NBIS)	Sweden	
45	Ryan	O'Connor	Independent	Independent	Ireland	
46	Ugochi	Okengwu	Faculty Member / Researcher	University of Port Harcourt	Nigeria	
47	Allison	Olsen	Digital Archivist	Children's Hospital of Philadelphia (CHOP)	United States	
48	Monica	Palmero Fernandez	Research Practice Coordinator	University of Oxford	United Kingdom	0000-0001-5164-9115
49	Rajeev	Pillai	Senior Project Manager	Beth Israel Deaconess Medical Center (BIDMC)	United States	
50	Nadège	Poncet	Agentic AI consultant	Witvio	Switzerland / France	0000-0002-0199-994X
51	Ronit	Purian	Researcher	Climatact	Israel	
52	Tovo	Rabemanantsoa	IT Manager	Institut National de Recherche pour l'Agriculture, l'Alimentation et l'Environnement (INRAE)	France	0000-0002-8362-9474
53	Trish	Radotic	RDA Community Manager	Research Data Alliance (RDA) / Australian Research Data Commons (ARDC)	Australia	
54	Leonore	Reiser	Biocuration Scientist	Phoenix Bioinformatics	United States	

55	Davide	Rizzo	Junior Professor in Landscape Agronomy	Institut de Recherche pour le Développement (IRD)	France	
56	Abiel	Roche-Lima	Professor	University of Puerto Rico System	United States	
57	Chokri	Ben Romdhane	IT Senior Analyst	Centre National Universitaire de Documentation Scientifique et Technique (CNUDST)	Tunisia	0000-0001-5437-6706
58	Marco	Rorro	AI Solutions Architect	EGI Foundation	Netherlands / Italy	0009-0002-6260-1007
59	Anna	Sackmann	Data Services Librarian	University of California, Berkeley (UC Berkeley)	United States	0000-0002-3852-6951
60	Habiba	Sarhan	Model Behaviour and Responsible AI Lead	GaiaLogic AG / Technical University of Munich (TUM)	Germany	
61	Bernd	Saurugger	Researcher	Technische Universität Wien (TU Wien)	Austria	
62	Hugh	Shanahan	Professor of Open Science	Royal Holloway, University of London	United Kingdom	
63	Jonathan	Starr	Executive Director	Open Source Endowment	United States	
64	Frankie	Stevens	Director	Research Infrastructure Services	Australia	
65	Joanne	Stocks	Senior Research Fellow	University of Nottingham	United Kingdom	0000-0002-7800-6002
66	Veronika	Stoka	Professor	Jožef Stefan Institute (JSI)	Slovenia	
67	Janine	Strandberg	Data Steward	Delft University of Technology	Netherlands	0000-0002-0336-5035
68	Jenifer	Tabita Ciuciu-Kiss	Data Engineer	Universidad Politécnica de Madrid (UPM)	Switzerland	
69	Annerose	Tartler-Ostrizek	Data Steward	FedOSC Belgium / Royal Belgian Institute of Natural Sciences (RBINS)	Belgium	
70	Andrew	Treloar	Independent	DRI Consultant	Australia	0000-0002-8911-3081
71	Slava	Tykhonov	Head of AI	Committee on Data of the International Science Council (CODATA)	France	0000-0001-9447-9830
72	Raymond	Uzwyszyn	Director, Research and Technology	University of California, Riverside (UCR)	United States	
73	Filippo	Vasone	Data Steward	University of Bologna (Unibo)	Italy	0009-0000-9125-8813
74	Alex	Wade	Consultant	DataCite	United States	0000-0002-9366-1507
75	Benjamin	Wright-Jones	Software Engineer	Microsoft	United Kingdom	
76	Ben	Wu	Principle Consultant	NetApp	Australia	
77	Qi	Zhang	Researcher	Research Organization of Information and Systems (ROIS)	Japan	
78	Olga	Zubova	Independent Consultant	Independent Consultant (Standards and Policy)	Peru	
79	Carlo Maria	Zwölf	Data Infrastructure Director	Laboratory for the Study of the Universe and eXtreme phenomena (LUX), Observatoire de Paris	France	0000-0002-5762-6747

Appendix D: Potential Data Director Blueprint Implementors

Based on discussions with working group members, there is interest in further developing this draft Blueprint and using it to set up implementations of the tool, or variations of it. Microsoft has also confirmed its intention to adopt the Blueprint within its own commercial context, bringing significant organisational capacity and reach to early implementation efforts.

Name/Affiliation	Description
Paty Buendia Lifetime Omics	A separate initiative associated with the EOSC-Future/RDA Artificial Intelligence and Data Visitation (AIDV) Working Group and the related TIGER project DV4RDA is developing software components focused on implementing data visitation policies in practice. This includes a tool called Vantage, with plans to extend to restricted access data. Members of the working group are also discussing potential use of this Blueprint within a new agentic AI working group whose scope of work is currently being developed. They have expressed willingness to collaborate with implementers.
Francis P. Crawley EOSC-Future/RDA Artificial Intelligence & Data Visitation Working Group (AIDV); RDA BIDT IG, Human Genome Project II (HGP2); AI4GHR; GCPA; SIDCER	The Data Director Blueprint has been adopted as a reference resource for the Global Ethics and Regulatory Preparedness for Emerging Health Technologies (GERP) initiative. It was implemented as part of the National Conference and Workshop: Towards a Safe, Equitable and Accountable Future of AI-Health Ecosystem in Indonesia, held at the Ministry of Health in Jakarta on 8–9 June 2026. The conference brought together ethics committees, regulators, and national health bodies to develop policy, SOPs, and assessment frameworks across AI and emerging health technologies. GERP will continue to reference the Blueprint in work across other countries, with a particular focus on the Global South.
Francis P. Crawley EOSC-Future/RDA Artificial Intelligence & Data Visitation Working Group (AIDV); GCPA; SIDCER	Members of the EOSC-Future/RDA AIDV Working Group are discussing potential use of this Blueprint within a new agentic AI working group whose Statement of Work is currently being developed. They have expressed willingness to collaborate with implementers.
Joe Heffer University of Sheffield Monica Palmero Fernandez University of Oxford	Two UK-based universities, Sheffield and Oxford, have expressed interest in exploring implementation of the Data Director. Both have indicated that collaboration between institutions could be beneficial. While no formal implementation plan has yet been established, they have indicated willingness to engage further with the Blueprint community.
Marco Rorro EGI Foundation	The EGI Foundation and its federated members provide research services ⁶⁸ that could serve as components for implementing the Data Director, including identity and secret management infrastructure. A federated large language model inference service is also expected to become available in the near term, which could provide the model runtime required for deployment.
Slava Tykhonov CODATA	A member of the Data Director Agentic AI Tool Blueprint Working Group is working on data management planning tools and have indicated they do not plan to implement the Data Director in full but have existing components that align with its requirements. These include an ODRL-based implementation addressing capability C12 and Croissant support addressing R2 and R4. They have expressed willingness to share these components and collaborate with other implementers.

⁶⁸ <https://www.egi.eu/services/research/>

<p>Raymond Uzwyshyn University of California</p>	<p>A member of the Data Director Agentic AI Tool Blueprint Working Group is developing a prototype that takes a persona-based approach as a complement to the capabilities-based model set out in this Blueprint. The prototype draws on the working group's outputs while adopting a different structural approach. Feedback from the wider community is being sought over the coming months and the developer has indicated willingness to collaborate with implementers.</p>
<p>Carlo Maria Zwölf LUX Observatoire de Paris</p>	<p>A suggestion has been raised to use the Blueprint specification itself as an informal test of its clarity and completeness, by providing it to a coding agent to assess whether it is sufficiently detailed to support AI-assisted implementation. This experiment has not yet been carried out but offers a practical means of evaluating the specification's rigour and identifying gaps ahead of broader implementation.</p>

Appendix E: Tool Usage

Working Group sessions were facilitated using [Zoom](#),⁶⁹ with collaborative notes taken in [Google Docs](#).⁷⁰ [Draw.io](#)⁷¹ and [Miro](#)⁷² were used for collaborative working and generating process flow and Blueprint architecture diagrams. Final process flow images were generated using [Claude.ai](#).⁷³

In line with the [RDA's Guidance on AI Tools Usage](#),⁷⁴ Claude Sonnet 4.5 (Pro) was used for data analysis and writing assistance throughout this Blueprint. The model was configured to ensure that input and output data are not used for model training. All AI-assisted analysis and AI-generated text have been reviewed, validated, and edited by the authors and working group members to ensure accuracy and completeness.

Appendix F: Partners and Document Information

This Blueprint is a collaborative effort between the Research Data Alliance and Microsoft, with contributions from participants of both organisations who did not receive any compensation for their involvement. The content generated from Working Group sessions has been revised from collaborative notes, transcripts and video recordings.

About the RDA

The [Research Data Alliance \(RDA\)](#) was launched as a community-driven initiative in 2013 with the vision that researchers and innovators can openly share and re-use data across technologies, disciplines, and countries to address the grand challenges of society. The RDA's mission is to build the social and technical bridges that enable that vision, accomplished

⁶⁹ <https://www.zoom.com/>

⁷⁰ <https://docs.google.com/>

⁷¹ <https://www.drawio.com/>

⁷² <https://miro.com/>

⁷³ <https://claude.ai/new>

⁷⁴ <https://www.rd-alliance.org/about/code-of-conduct/rda-guidance-on-ai-tools-usage/>



through the creation, adoption and use of the social, organisational, and technical infrastructure needed to reduce barriers to data sharing and exchange.

As of June 2026, the RDA comprises a 16,700+ member-strong community of researchers, data professionals, publishers, funders and policymakers that collaborate in working groups, interest groups and communities of practice to create recommendations and outputs.

[Individual membership](#)⁷⁵ is free of charge and open to all who share the RDA's Guiding Principles. To get involved at the organisational level, explore our [organisational and affiliate membership options](#).⁷⁶

About Microsoft

[Microsoft Corporation](#)⁷⁷ is a multinational American technology company recognised for shaping the evolution of personal and enterprise computing. Founded in 1975 and headquartered in Redmond, Washington, the company initially revolutionised software accessibility through its early operating systems. Over the decades, Microsoft expanded its portfolio to encompass a broad spectrum of technologies, including productivity, software, cloud infrastructure, gaming, Artificial Intelligence and Quantum Computing. With a longstanding emphasis on innovation and digital transformation, Microsoft continues to play a pivotal role in defining the future of the tech industry.

Licence and Contact

Licensing: All materials are licensed under Creative Commons Attribution 4.0 International (CC BY 4.0), allowing reuse with appropriate citation of this report.

Contact: For questions about this Blueprint, the RDA Working Group and facilitation programme, please contact the RDA Secretariat at [secretariat@rda-foundation.org].

⁷⁵ <https://www.rd-alliance.org/membership/individual-membership/>

⁷⁶ <https://www.rd-alliance.org/membership/organisational-membership/>

⁷⁷ <https://www.microsoft.com/>