

Security in Aridhia's Digital Research Environment



Author: **Robert Bryce**
Date: **14 September 2021**
Version: **Draft2**

Table of Contents

PREFACE	2
THE EVOLUTION OF RESEARCH ENVIRONMENTS	2
<i>Secure data platforms and safe havens</i>	2
<i>Meeting the needs of Health data scientists</i>	2
<i>Increasing information security requirements</i>	3
<i>Acceptance of cloud services</i>	4
<i>Emergence of new models</i>	4
ARIDHIA'S DIGITAL RESEARCH ENVIRONMENT	4
Aims	4
DRE security controls and the Five Safes framework	5
<i>Safe People</i>	5
<i>Safe Data</i>	6
<i>Safe Projects</i>	6
<i>Safe Settings</i>	6
<i>Safe Outputs</i>	7
SUMMARY	7

Preface

As the Aridhia Digital Research Environment (DRE) has developed from its early days as an on-premises implementation through to a modern cloud-native service, we have tried to balance two principal requirements from our user community. That is, to meet the functionality required by data scientists to conduct their research whilst also satisfying the governance and controls required by information security teams within a healthcare setting.

Data scientists want to creatively analyse varied types of data, often at scale and at speed, using a variety of paradigms and tools. They work in what is essentially a research and development role, therefore freedom to operate on some level is required to produce outcomes which can be translational in their field of study. Information security teams however, must consider the risk to privacy and security of any sensitive data, given the organisation's obligations to individuals, data contributors, and the legal authorities.

In this paper we will discuss the design and implementation choices made in Aridhia's Digital Research Environment to achieve this balance of usability and security.

The evolution of research environments

SECURE DATA PLATFORMS AND SAFE HAVENS

Early versions of research environments were created in various settings, with technology implementations which varied from being relatively open (in terms of security controls) to heavily locked down. In clinical/healthcare scenarios they tended towards being heavily locked down (often called "Safe Havens") and were characterised as follows:

- Project and users would be vetted prior to access to the environment.
- On-premises deployments within an organisation's own data centre, managed by an internal IT department.
- Infrastructure typically consisted of virtualised servers which users would access through a remote desktop connection.
- Users would then see a locked down virtual desktop with a selection of databases and analytical tools.
- In some cases, resources could only be accessed from certain physical locations and equipment i.e., users would enter a secure room and use a specific desktop device.
- Analytical output would be vetted prior to release outside the environment.

While these implementations satisfied the needs of many research scenarios, others required more flexibility which the heavily locked down safe haven model struggled to meet. This, combined with developments in health data science, information security, and in cloud services, meant that other implementations of a research environment were required.

MEETING THE NEEDS OF HEALTH DATA SCIENTISTS

Data science in a healthcare context brings together different disciplines which in turn need to be reflected in the research environment. Data science as a discipline has an emphasis on an open and collaborative culture, with the use of open-source tools and programming languages which allow code, packages, and techniques to be shared and reproduced within the wider community. Data scientists are often highly proficient in technology and programming and expect a degree of leeway in configuring an environment to match their preferred ways of working.

Researchers from a more traditional healthcare background will bring subject matter expertise in their clinical field and most likely, statistical expertise too. Their technical and programming knowledge may be more limited than someone from a pure data science background, however they will share the desire to collaborate safely within a wider community.

A lead researcher/data scientist has therefore had to consider:

- How can I prepare and analyse data using a variety of programming languages and statistical/analytical packages, according to the project needs and researcher skill sets?
- How can I bring my own software tools or analytical pipelines to the environment?
- Can I invite colleagues from other institutions to work with me in the research environment?
- Can I bring my own data to the environment e.g., as captured in an electronic case report form (eCRF) or similar tool during a clinical research project?
- Can we pool data across environments and geographies? Data for certain health conditions can be limited due to the size of the patient cohort in one geography, hence the need to pool data to create a meaningful study size.
- How can I make output research datasets confirm to the [FAIR data principles](#) (that is to make them Findable, Accessible, Interoperable, Re-usable) for future benefit?
- How can I take advantage of machine learning and emerging healthcare technology standards such as [FHIR](#)?

Historically it has been difficult to reconcile many of these requirements in a traditional safe haven, which meant that some projects gained little traction.

INCREASING INFORMATION SECURITY REQUIREMENTS

As the number of security and privacy threats has grown, so has public awareness and concern. To address this, technology, standards and accreditations have been developed as well as legislation such as the General Data Protection Regulation (GDPR) which came into effect in 2018. Within research environments in particular, we have also seen the emergence of the [Five Safes framework](#), developed in the UK in 2017.



The ISO 27001 standard is now well accepted as the recognised baseline for establishing an information security management system. An organisation is accredited by an independent third party in order to receive the certification, providing assurance that the standard is being met. The ISO 27701 extension applies additional privacy management controls to 27001, should an organisation determine that this is required. Within the US, processing of electronic personal health information requires compliance with the Health Insurance Portability and Accountability Act (HIPAA). There is no formal certification process for HIPAA, however the HITRUST certification has been developed which aims to rationalise several internationally accepted security and privacy-related regulations, standards, and frameworks (ISO, NIST, PCI, HIPAA, and COBIT) into one certification.

GDPR sets out legal obligations for processors and controllers of personal data, which centre around:

- Lawfulness, fairness, and transparency
- Purpose limitation
- Data minimisation
- Accuracy
- Storage limitation
- Integrity and confidentiality (security)
- Accountability.

These obligations will not be explored in detail here; however, it is important to state that since coming into effect in 2018, organisations must be in a position to comply with this legislation and if necessary, demonstrate this to the

relevant authority within their jurisdiction. For example, a breach of personal data must be reported within 72 hours to the supervisory authority. Within the EU, non-compliance to GDPR can result in fines up to €20 million or 4% of annual global turnover.

Whilst not an accredited standard in the manner of ISO 27001 or HITRUST, the Five Safes framework has emerged as a recognised way of approaching safe research and data sharing. The framework was developed by the UK Data Service with the ONS and HMRC to provide safe research access to data. The framework can be summarised as “SAFE People accessing SAFE Data within a SAFE Settings to undertake SAFE Projects resulting in SAFE Outputs.” It has been adopted in the UK and elsewhere by organisations including [Health Data Research UK](#) (HDR-UK) and [National Institute for Health Research Design Service](#) (NIHR).

ACCEPTANCE OF CLOUD SERVICES

Early research environments were usually built on-premises or using specialised cloud vendors which served a specific territory and customer base. As the hyperscale cloud vendors like AWS and Microsoft established data centres in more countries and achieved the necessary certifications (in quality, information security, data centre, environmental, etc.) organisations such as the NHS in the UK approved their use for the processing of healthcare data.

Cloud vendors in turn invested in technologies which the healthcare and research communities could take advantage of, such as machine learning and Fast Healthcare Interoperability Resources (FHIR) services.

Research environments could also take advantage of other common benefits of the cloud such as auto-scaling, pay-as-you-go, and significant capacity and variety across storage, compute, and other aspects of an infrastructure. Significant investment is required to provide high-performance computing services within a highly secure environment, which many organisations would struggle to fund in an on-premises deployment therefore the ability to quickly scale-up and then scale-down can offer a much more economical model.

Whilst there will be organisations which will continue to require on-premises for certain use cases, the cloud has now been accepted and in many cases adopted in the healthcare research field.

EMERGENCE OF NEW MODELS

Given the background we have discussed, the next generation of research environments are emerging which build on earlier models, whether developed in public, industry, or philanthropic settings. The concept of [Trusted Research Environments](#) (TREs), particularly within the UK, is often used to categorise such environments.

TREs share many of the same characteristics as earlier environments, with an emphasis on consistently and transparently applying the Five Safes framework, such that public concerns around how data is being used can be addressed. TREs have no formal accreditation process albeit this is a topic for discussion within the UK community.

Aridhia's Digital Research Environment

Aims

In bringing our Digital Research Environment to market, Aridhia has focussed on providing a service which meets the following requirements:

- That it follows the research lifecycle, from finding and accessing data through our *FAIR* service, through to collaboratively analysing and sharing results in our *Workspaces* service
- That it is aligned to the various roles within the research lifecycle (i.e., Principal Investigators and Study leads, Data Stewards and Researchers) who should be able to operate in a largely self-serve manner

- That Aridhia be a data processor in providing the service and not a data controller (which remains the responsibility of the customer/user). Aridhia has no rights over any data hosted within the DRE, other than service usage data required for support (e.g., number of concurrent users).
- That the service be designed for collaboration between researchers and data scientists with mixed technical skills, often working together in a multi-institutional setting. The DRE is therefore hosted on a public cloud (Microsoft Azure), internet connected (subject to additional controls which can be applied) and accessed over a standard web browser
- That the service is for the analysis of anonymised and pseudonymised data (which for the purposes of GDPR is still treated as personal data). Should customers wish to analyse identifiable data, then additional security controls would be applied
- That it be operated in accordance with recognised information security certifications and applicable laws such as GDPR
- That customer organisations accept a shared responsibility model with Aridhia. This means that customers will accredit their researchers and ensure that users follow the terms and conditions of the service, with Aridhia taking responsibility for security controls and measures in the service (e.g., virus scanning of uploads, two-factor authentication, secure workspace boundaries, vulnerability scanning of containers, etc.).

It should also be emphasized that Aridhia is providing a managed service to its customers. Historically, many organisations have engaged their internal IT team to build a research environment. This can be difficult to sustain in the long-term however, given the significant level of investment and effort required. As the DRE is a managed service, Aridhia performs all deployments, upgrades, OS patching, service monitoring and provides end user support, as well as providing regular feature and functionality upgrades.

DRE security controls and the Five Safes framework

We will now go on to discuss the various information security features which are embedded into the DRE to meet the requirements we outlined in the previous section. We will use the Five Safes framework as a reference point for this, as while it is not a formal accreditation or agreed list of controls, it is useful as a means of categorising the different risks which exist and of the measures which are being taken to manage them.



SAFE PEOPLE

Users of the DRE will self-register and then follow a two-factor authentication process to login. As part of the authentication process, users accept the terms and conditions for the environment, which includes abiding by the local data protection legislation and not attempting to bypass the built-in security controls.

Customer organisations will also have their own accreditation process, to ensure that users allowed to access data and projects within the DRE are aware of their responsibilities and have been given suitable training. Aridhia works with customers in this on-boarding and training activity, however it is the responsibility of the customer organisation to ensure that users are accredited/approved to access their DRE.

All user access and subsequent activities in the DRE are captured in audit logs, which can be viewed by administrators within the relevant organisation to ensure that practices are being followed correctly.

As mentioned earlier, customer organisations and their users must accept a shared responsibility model, where they follow the terms and conditions of using the DRE. This includes not attempting to bypass or override the built-in

security controls. Not all attempts to override these controls are deliberate of course, therefore training and awareness is required to help guide users appropriately.

SAFE DATA

Organisations will determine which datasets are suitable for storage and analysis within the DRE, with Data Sharing Agreements in place between Aridhia (acting as the Data Processor) and organisations (acting as the Data Controller).

Datasets can be held in Aridhia's FAIR service and are released into Workspaces for analysis following a Data Access Request and approval process.

Data ingress into the DRE can be done via API (using a secure token) or through the UI, with all files going through a malware scanning process. In the Workspaces service, files which have been malware scanned then go through an inbound airlock approval process where they can be approved or rejected for entry into a specific project workspace.

From the FAIR service, personal identifiable information (PII) can first be pseudonymised or synthesised before being passed into a workspace for analysis. Data within the Workspaces service will therefore be anonymised or pseudonymised. Should a customer want to hold PII within a workspace, additional controls can be implemented including locking down access to a specific IP address range. It should be noted that DRE environments are single tenanted i.e., there is no shared infrastructure between different customers, their respective datasets and other resources.

SAFE PROJECTS

Projects will be reviewed by the customer organisation for scientific merit and benefit, before being approved. Once a project is approved, the role-based access control mechanism within the DRE is used to determine how projects and related resources are controlled and managed:

- A Principal Investigator or Lead Researcher will be assigned as a workspace administrator for their project, determining which other researchers to invite, what data to bring in, and what, if anything, can be exported from the workspace.
- A researcher will be assigned a *standard* user role while users who will just be required to upload data can be given a *contributor* role.
- At an organisational level, a *tenant administrator* is assigned, who will determine how workspaces will be allocated to projects which have been approved. The role is often fulfilled by a data steward team who manage the research environment on behalf of the organisation.
- A similar set of roles apply to datasets within the FAIR service.

SAFE SETTINGS

Certification and standards

Aridhia's DRE is a certified environment and therefore independently audited each year to the ISO 27001 information security standard.

Aridhia is also soon to acquire HITRUST accreditation which we expect to receive in Q4 2021. Our services are penetration tested several times each year by independent third parties.

The service is built to a documented software development lifecycle, with secure coding guidelines that follow the [OWASP top-10](#) guidelines.

The DRE is deployed to Microsoft Azure regions as chosen by the customer organisation, given their preferences for where data and resources should be geographically located. Azure regions are highly controlled environments which are certified to numerous information security, quality, and cloud service standards.

Security boundaries

Within the DRE, boundaries are enforced through network security groups and the role-based access permissions model mentioned earlier. It should be noted that if users are members of multiple workspaces, for example *abc* and *def*, they can only access resources for *abc* from within that particular workspace.

As discussed earlier, researchers want flexibility within the environment. For example, they may want additional software tools hosted on a virtual machine running in a workspace. This is controlled through configuration of an allow-list, which determines which internet resources (if any) can be accessed from a workspace to allow a software download (which is then virus scanned). Containerised applications can also be brought to the DRE, which are scanned for vulnerabilities before being made accessible to end users.

Encryption

All user access is via HTTPS URL protected by a rooted certificate issued by DigCert SHA2 Secure Server CA, utilising sha256RSA signature algorithm with sha256 signature hashing algorithm. TLS 1.2 protocols or above.

All internal network traffic is protected by HTTPS or TLS 1.2 or above protocols.

Encryption at rest using FIPS 140-2 compliant 256 AES encryption for storage accounts and virtual machine disks.

Managed service

The DRE is a managed service, whereby Aridhia performs the following operations:

- OS patching (scheduled and ad hoc in the event of emergency updates)
- Configuration and monitoring of intrusion prevention and detection systems
- Nightly back-ups of the environment
- All support teams have separate privileged admin accounts and these require 2FA. All support team actions are logged
- Regular audits of the privileged accounts
- All Aridhia employees who have access to the operational environment go through a criminal record check
- Mandatory security training at inductions and periodic refresher training for all employees
- All changes to the platform are subject to change control
- Incident management and CSIRT processes
- Monthly audit of key ISO27001 controls
- Quarterly BCP/DR exercises.

SAFE OUTPUTS

Researchers must use the supported data export mechanism, the outbound airlock, to extract data, results, code etc. from the DRE. Workspace administrators (i.e., the PI/research lead for the project) must approve outbound airlock requests before any download can be completed. Note that the airlock mechanism is also used as a safe control for transfers between workspaces, as well as outside of the DRE.

The customer organisation is therefore responsible for output checking, given their role as the data controller.

Other mechanisms to export data such as copy/paste from a browser are not supported as part of the terms and conditions of the DRE. Aridhia has chosen not to disable copy/paste from the browser as this would degrade the overall user experience. Being unable to copy and paste would be very frustrating for a data scientist trying to write R or Python code for example. Instead, warning messages are sent to the user reminding them not to copy outside of the DRE and additional logging is in place to record these actions.

Summary

In this paper we have explored how research environments have evolved and the factors influencing their development, particularly in data science, security and the cloud.

The Aridhia DRE has been built with a “security by design” approach, whereby we offer features for researchers and health data scientists which have been developed within a secure foundation. There are many layers to a robust security design and in the context of the DRE, we have framed these layers using the Five Safes framework as a reference point.

Security is a broad and continuously evolving topic and at Aridhia we are now considering the next set of controls to add to the existing foundation we have in place, such as further certifications, adding synthetic data to the DRE, and how to evolve security in a federated analysis model.