



Beyond One Million Genomes

## D4.1 - Secure cross-border data access roadmap - 1v0

Project Title (grant agreement No)	Beyond One Million Genomes (B1MG) Grant Agreement 951724		
Project Acronym	B1MG		
WP No & Title	WP4 - Federated Secure Cross-border Technical Infrastructure		
WP Leaders	Tommi Nyrönen (CSC), Ilkka Lappalainen (CSC), Bengt Persson (UU), Sergi Beltran (CNAG-CRG)		
Deliverable Lead Beneficiary	4 - CSC		
Deliverable	D4.1 - Secure cross-border data access roadmap - 1v0		
Contractual delivery date	31/05/2021	Actual delivery date	28/05/2021
Delayed	No		
Authors	Dylan Spalding (CSC), Tommi Nyrönen (CSC), Heikki Leväslaiho (CSC), Riku Riski (CSC), Ilkka Lappalainen (CSC), Bengt Persson (UU), Sergi Beltran (CNAG-CRG)		
Contributors	Regina Becker (UNILU), Juan Arenas (ELIXIR Hub)		
Acknowledgements (not grant participants)			
Deliverable type	Report		
Dissemination level	Public		



Beyond One Million Genomes

B1MG has received funding from the European Union's Horizon 2020 Research and Innovation programme under grant agreement No 951724



## Document History

Date	Mvm	Who	Description
01/04/2021	0v1	Dylan Spalding (CSC)	[Initial draft circulated to WP participants for feedback]
20/05/2021	0v2	Nikki Coutts (ELIXIR Hub)	Version circulated to B1MG-OG, B1MG-GB & Stakeholders for feedback
28/05/2021	0v3	Dylan Spalding (CSC)	[B1MG-OG/GB/stakeholder comments addressed]
28/05/2021	1v0	Juan Arenas (ELIXIR Hub)	Final version uploaded to the EC Portal

## Table of Contents

- [1. Executive Summary](#)
- [2. Contribution towards project objectives](#)
- [3. Methods](#)
- [4. Description of work accomplished](#)
  - [4.1 \[Section headers\]](#)
    - [4.1.1 \[Sub-section header\]](#)
- [5. Results](#)
- [6. Discussion](#)
- [7. Conclusions](#)
- [8. Next steps](#)
- [9. Impact](#)
- [10. References](#)

## 1. Executive Summary

A roadmap is required to communicate the plans for development of the infrastructure to provide secure cross border data access to genetic and phenotypic data to interested stakeholders. The roadmap is a live document, with the initial deliverable being version1, which will evolve as the requirements of more use cases and updates to technology and standards are taken into account.

The main aim of the roadmap presented here is to define the timelines to provide a proof of concept demonstrator of secure cross border data access, describe the technologies required to



develop such a demonstrator, and to detail the work that has gone into determining the roadmap and requirement for the proof of concept demonstrator. The roadmap also outlines how it is expected that the work carried out towards the initial proof of concept will support the developing roadmap for 2022 and beyond, as the proof of concept evolves from a demonstrator with synthetic data to a production system accessing real genetic and phenotypic data.

## 2. Contribution towards project objectives

With this deliverable, the project has reached or the deliverable has contributed to the following objectives/key results:

[Select 'Yes' (at least one) if the deliverable contributed to the key result, otherwise select 'No'.]

Objective/Key Result No & Description	Contributed
<b>Objective 1:</b> Engage local, regional, national and European stakeholders to define the requirements for cross-border access to genomics and personalised medicine data	
1: B1MG assembles key local, national, European and global actors in the field of Personalised Medicine within a B1MG Stakeholder Coordination Group (WP1) by M6.	No
2: B1MG drives broad engagement around European access to personalised medicine data via the B1MG Stakeholder Coordination Portal (WP1) following the B1MG Communication Strategy (WP6) by M12.	No
3: B1MG establishes awareness and dialogue with a broad set of societal actors via a continuously monitored and refined communications strategy (WP1, WP6) by M12, M18, M24 & M30.	No
4: The open B1MG Summit (M18) engages and ensures that the views of all relevant stakeholders are captured in B1MG requirements and guidelines (WP1, WP6).	No
<b>Objective 2:</b> Translate requirements for data quality, standards, technical infrastructure, and ELSI into technical specifications and implementation guidelines that captures European best practice	
<b>Legal &amp; Ethical Key Results</b>	
1: Establish relevant best practice in ethics of cross-border access to genome and phenotypic data (WP2) by M36	Yes
2: Analysis of legal framework and development of common minimum standard (WP2) by M36.	Yes
3: Cross-border Data Access and Use Governance Toolkit Framework (WP2) by M36.	Yes
<b>Technical Key Results</b>	



4: Quality metrics for sequencing (WP3) by M12.	No
5: Best practices for Next Generation Sequencing (WP3) by M24.	No
6: Phenotypic and clinical metadata framework (WP3) by M12, M24 & M36.	No
7: Best practices in sharing and linking phenotypic and genetic data (WP3) by M12 & M24.	No
8: Data analysis challenge (WP3) by M36.	No
<b>Infrastructure Key Results</b>	
9: Secure cross-border data access roadmap (WP4) by M12 & M36.	Yes
10: Secure cross-border data access demonstrator (WP4) by M24.	Yes
<b>Objective 3:</b> Drive adoption and support long-term operation by organisations at local, regional, national and European level by providing guidance on phased development (via the B1MG maturity level model), and a methodology for economic evaluation	
1: The B1MG maturity level model ( WP5) by M24.	No
2: Roadmap and guidance tools for countries for effective implementation of Personalised Medicine (WP5) by M36.	No
3: Economic evaluation models for Personalised Medicine and case studies (WP5) by M30.	No
4: Guidance principles for national mirror groups and cross-border Personalised Medicine governance (WP6) by M30.	No
5: Long-term sustainability design and funding routes for cross-border Personalised Medicine delivery (WP6) by M34.	No

### 3. Methods

This document describes the evolving roadmap to provide the infrastructure required to support secure cross border access to the genetic, clinical, phenotypic and pedigree data generated within the European Union by the end of 2022. The document defines the current standards, functionalities, and components that are required for the Proof of Concept (PoC) demonstrator required in 2022. It outlines the status of each of these, and whether further development is required to ensure they can form a functional part of the proposed infrastructure for the PoC demonstrator.

To ensure the infrastructure solution proposed by WP4 supports all the requirement of other work packages in both Beyond 1 Million genomes (B1MG) and the working groups of the 1+ Million Genomes (1+MG), a scoping paper<sup>1</sup> was written and three workshops were arranged

<sup>1</sup>

<https://docs.google.com/document/d/1L5imuKcL0wZNQQ1vQmPAXpR8SG42EL4iq1HZJH87Yv0/edit>



(one on the requirements for synthetic data<sup>2</sup>, one on the rare disease use-case<sup>3</sup>, and one on Ethical, Legal, and Societal Issues (ELSI) and Data Protection by Design with WP2<sup>4</sup>).

The synthetic data workshop was arranged in December 2020 to investigate the availability of synthetic datasets and the requirements of the different 1+MG use-case working groups for synthetic datasets. By matching the synthetic datasets to the requirements of the use cases, and the capability of each use-case to utilise the proposed infrastructure in time for the PoC demonstrator, a synthetic dataset was chosen to be used along with the rare disease use case.

The scoping paper was written by 1+MG working group 5 and B1MG WP4. The paper defines the five core functionalities required to support secure cross border data access. The paper listed the early adopter nodes that can deploy the proposed infrastructure for the PoC demonstrator, and proposed the first use-case to be addressed by the demonstrator to be the rare disease use case.

A rare disease workshop was held, where the PoC was described, and gap analysis performed to identify gaps where the PoC would not support the proposed rare disease use case. To ensure maximum alignment with existing rare disease projects, such as the European Joint Program on Rare Disease<sup>5</sup> (EJP-RD which is a driver project for the Global Alliance for Genomics and Health<sup>6</sup> (GA4GH), ensuring the PoC utilised relevant standards applicable to rare disease was agreed.

The joint workshop with WP2 on ELSI and Data Protection by Design detailed the proposed Proof of Concept infrastructure, the nine main principles of ELSI that must be satisfied, and outlined the process to ensure that the proposed infrastructure conforms to those principles and General Data Protection Regulation (GDPR).

## 4. Description of work accomplished

In this section, the technological solutions proposed are outlined, including the standards, functionalities, and components chosen to provide the core infrastructure. The core components are the separate components that currently exist and are being actively developed and maintained through a variety of other sources, and these components can be comprised of one or more services or applications. These components together provide the five functionalities that were defined in the scoping paper as being required to provide the required infrastructure. Linking the services and components together is done via global standards where possible, ensuring that the infrastructure is interoperable with external infrastructures where possible.

---

<sup>2</sup> <https://docs.google.com/document/d/1PVUDsVccjuLHynSX5fG1aDB8dCHiga9woGa7kvF5bE8/edit>

<sup>3</sup> [https://docs.google.com/document/d/1T4x575AgU3QPzBHS5wYb20ANasII\\_8\\_mSJ-Hd4WrPAA/edit](https://docs.google.com/document/d/1T4x575AgU3QPzBHS5wYb20ANasII_8_mSJ-Hd4WrPAA/edit)

<sup>4</sup> <https://docs.google.com/document/d/1kCiCftFTIODjRZOZ-MTktCCCqdgohZDebTEGIIAg588/edit>

<sup>5</sup> <https://www.ejprarediseases.org/>

<sup>6</sup> <https://www.ga4gh.org/>

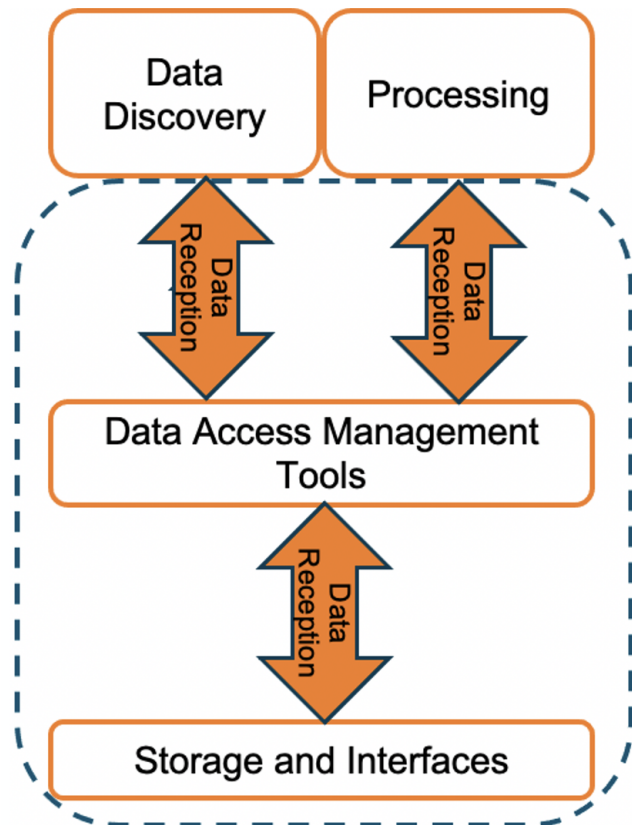


## 4.1 Infrastructure Functionalities

To provide the two main user-facing functionalities, the Data Discovery and Processing functions, the three other functions are required – Storage and Interfaces, Data Access management Tools, and Data Reception.

As described in the scoping paper, the data processing functionality can include running analytical pipelines (see later under GA4GH), but can also include visualisation or other processing activities. For example, the RD-Connect Genome-phenome Analysis Platform (GPAP) can be regarded as a processing activity as it manipulates the data it accesses to allow the user to visualise or analyse these data. As detailed in the scoping paper, the infrastructure responsibilities of WP4 are the functionalities within the dashed line in Figure 1, which can in general be seen as a Federated EGA instance with extended functionalities (see Table 1 for functionalities of a FEPA instance), such as the Resource Entitlement Management System<sup>7</sup> (REMS).

For this proof of concept, it is put forward that the infrastructure proposed by B1MG demonstrates both user-facing functionalities, supported by the core infrastructure in a federated environment.



**Figure 1:** Functionalities described in the scoping paper, and their relationship to the user facing functionalities of Data Discovery and Processing

<sup>7</sup> <https://github.com/CSCfi/remis>

## 4.2 Main Services

### 4.2.1 Beacon

The Beacon service provides the Data Discovery functionality. It is a GA4GH standard and a service provided by the GPAP. It allows a user to query a controlled access resource without having gained prior approval from a Data Access Committee (DAC) to access that resource while respecting the privacy of the individuals whose data exists within that resource. As such, the Beacon currently allows queries based on the presence of a particular allele at a certain genomic locus. The Beacon may return a binary yes / no answer, an allele frequency of that particular allele within the queries resource, and additionally the dataset within the resource that the allele (or alleles) exist. The standard supports three levels of access, anonymous, registered, and controlled. Anonymous access is effectively public access, where a user needs to have an identity associated with their query. Registered access requires that the user has an identity, such as an ELIXIR identity, and controlled access requires that the user has an identity and has also been granted permission to query the particular Beacon by the data controller or data access committee. Depending on the individual Beacon, registered access may require that a 'bona fide' researcher attribute is set in the user's identity, which ensures the user is a member of a recognised institution. Each Beacon may determine the type of response for each access level, depending on the requirements of the DAC or data controller of the information within the Beacon.

### 4.2.2 Federated EGA

The Federated EGA (FEGA) in this context is a technical infrastructure that provides a solution for archiving and distributing genetic and phenotypic data to authorised users only. As such it has components, the archive or file storage, the user permissions service, and the data distribution service. The archive utilises the file storage infrastructure provided by the node, and ensures the files archived at the FEGA are encrypted using the Crypt4GH encryption standard so that the files can only be read by authorised users. The user permissions service stores the access permissions for each user who can access data held at the FEGA node. Each user has an access account, and the files or datasets that the user can access are associated with the user account. This account can be linked to an ELIXIR identity using the GA4GH passport standard. The data distribution service is the service that is responsible for distributing the data from the archive to the user or other authorised location or service. The data distribution service utilises the GA4GH htsget standard, which ensures the data is transmitted over https, a secure encrypted data transfer protocol. The FEGA provides the Data reception and Storage and Interfaces functionalities, and as part of these functionalities maintains logs of user access permissions, file location and access, and data distribution processes.



### 4.2.3 Federated GPAP

The federated Genome-phenome Analysis Platform allows users to interact with, and visualise genetic and phenotypic data and run analyses. For example, the user can search and visualise variants that occur in individuals with a certain phenotype or disease. The GPAP provides a Beacon service for data held within the GPAP. The GPAP also has a user permissions database that stores the access permissions. The access permissions service of the GPAP is currently being modified to ensure it is compatible with the GA4GH Passport standard. The GPAP provides the processing functionality.

### 4.2.4 Resource Entitlement Management System

The Resource Entitlement Management Systems (REMS) is a service that allows users to apply for data access, and data controllers or DACs to administer data access for data held in one or more FEAGA nodes. For example, after discovering a dataset of interest via a Beacon, a user is directed to make an access request to a file or datasets which is listed in REMS, and REMS then automatically notifies the DAC that a request has been made. The DAC can then log into REMS, and grant or deny the access request. The process of applying for dataset access can be configured by the DAC to support different procedures for applying for data access, for example ensuring the Data Access Agreement is signed. If the access request is granted, REMS updates the user permissions service at the appropriate FEAGA node(s), and notifies the users that they now have access to these data. REMS provides the Data Access Management Tools functionality.

## 4.3 GA4GH Standards Used in the Proof of Concept

In this roadmap, we propose to base the infrastructure technologies on global standards where appropriate as this ensures maximum interoperability both within Europe and worldwide. The Global Alliance for Genomics and Health (GA4GH) is an organisation that sets the technical standards and frames the policy for the responsible sharing of human genomic data within a human rights framework. Every GA4GH standard must go through a defined approval process (Appendix 1), which helps ensure that these standards address a real use-case and meet minimum security and ELSI standards as defined by the Data Security and Regulatory and Ethics foundational workstreams. GA4GH standards are also semantically versioned, which ensures that interoperability within the infrastructure can be maintained and tracked as the standards and technology evolves, ensuring sustainability. As defined in the scoping paper, there are 5 main functionalities proposed for the B1MG infrastructure which can be mapped to one or more existing GA4GH standards (table 1). An overview of each respective GA4GH standard is given below.

Functionality	GA4GH Standard	PoC Component	PoC Service
---------------	----------------	---------------	-------------





Data Discoverability	Beacon, DUO	REMS	Beacon
Data Reception	htsget, phenopackets	Data API	FEGA
Storage and Interfaces	VCF, SAM / BAM / CRAM, Crypt4GH, phenopackets		
Data Access Management Tools	AAI / Passports, DUO	REMS	ELIXIR AAI
Processing	TES, WES	IGV	GPAP

**Table 1:** Example mapping of GA4GH standards used in the Proof of Concept services and components.

A majority of the proposed components or services required for the PoC are open source (Table 2), with the exception of the GPAP.

Service or component	Open source	License	Link
Beacon	Yes	Apache 2.0	<a href="https://github.com/ga4gh-beacon/beacon-elixir/blob/master/LICENSE">https://github.com/ga4gh-beacon/beacon-elixir/blob/master/LICENSE</a>
Federated EGA	Yes	Apache 2.0	<a href="https://github.com/EGA-archive/LocalEGA/blob/master/LICENSE">https://github.com/EGA-archive/LocalEGA/blob/master/LICENSE</a>
Data API	Yes	Apache 2.0	<a href="https://github.com/EGA-archive/ega-data-api/blob/master/LICENSE">https://github.com/EGA-archive/ega-data-api/blob/master/LICENSE</a>
REMS	Yes	MIT	<a href="https://github.com/CSCfi/remis/blob/master/LICENSE">https://github.com/CSCfi/remis/blob/master/LICENSE</a>
IGV	Yes	MIT	<a href="https://github.com/igvteam/igv/blob/master/license.txt">https://github.com/igvteam/igv/blob/master/license.txt</a>
GPAP	No		
ELIXIR AAI	Yes	Apache 2.0	<a href="https://github.com/elixir-cloud-aa/elixir-cloud-aa/blob/dev/LICENSE">https://github.com/elixir-cloud-aa/elixir-cloud-aa/blob/dev/LICENSE</a>

**Table 2:** Licenses applying to the service or component implementation for the PoC.

### 4.3.1 Beacon

The GA4GH Beacon standard is a data discoverability standard that underpins the Beacon service described above. The current Beacon standard is relatively limited in its utility, allowing users just to query for a particular allele at a particular locus. Within the rare disease use case this allows users to determine if a node contains an individual who has the same variant as to one of interest, for example a candidate causative monogenic variant for a certain disease, which would allow the user to request full access to the underlying data. But as by definition rare disease variants make the individual more susceptible to re-identification, the use of the current Beacon standard within the rare disease use case is limited.

To try and address this, the new Beacon V2 standard is being developed. This major version release, which will need to be reviewed and approved by the GA4GH, will extend the Beacon to allow gene or region level queries, and phenotype queries. This significantly extends the utility of the Beacon for the rare disease use case to allow researchers to query



for nodes containing individuals who share variants in the same region or gene, share similar phenotypes or disease, or datasets that correspond to certain data use conditions.

### 4.3.2 Data Use Ontology

The Data Use Ontology (DUO) is an ontology that describes the data use conditions that are applied to one or more datasets held within a node. As the DUO is an ontology, it is by design machine-readable. The conditions can affect the type of secondary use that can be performed on the original data. For example, the data can be restricted to a certain location, or to research within a certain disease. Each release of the DUO is versioned, ensuring that historic DUO attributes are maintained when the definition of the term is modified to reflect changes in the legal and technological landscape around genomic data.

The DUO can be seen as part of the data Discoverability functionality, and the Data Access management Tools functionality. This is because the machine readable attribute of the DUO allows prospective users to query data sets based on the data use conditions - for example user A may wish to research disease X, so can filter out from available datasets all but that disease. The DUO can also be used to determine data access, for example by determining if the researcher who is applying for data access intends to use the data for a purpose allowed by the Data Access Committee or data controller.

### 4.3.3 htsget

htsget is a data streaming protocol that allows genetic data to be securely streamed from one location to another. It supports range-based queries, for example a user can request a single region, exon or gene to be streamed to a location. htsget is a secure protocol, with all data transmission running over https. The use of htsget can be used to restrict access to genomic regions, preventing users streaming data from certain genes for example.

### 4.3.4 Phenopackets

The Phenopacket standard is designed to allow phenotypic information to be exchanged in a consistent standardised and structured format. It links phenotypic information with individual, genetic, and disease information. Phenopackets are already used by the rare disease community, for example in the Solve-RD<sup>8</sup> project and RD-Connect<sup>9</sup>, which makes them ideal as a standard to drive the exchange of phenotypic information between nodes or services.

### 4.3.5 VCF, BAM / CRAM / SAM

The VCF and SAM / BAM / CRAM file formats<sup>10</sup> are used to store genetic data in a standard format. SAM / BAM / CRAM files store the complete genetic data for a genomic region,

---

<sup>8</sup> <http://solve-rd.eu/>

<sup>9</sup> <https://rd-connect.eu/>

<sup>10</sup> <https://samtools.github.io/hts-specs/>



including the depth of coverage for aligned data across these regions, and the quality of the sequenced reads. BAM is a binary version of a SAM file, while CRAM is a newer file format that compresses the data to a greater extent than SAM or BAM, reducing the amount of storage required for a set amount of genomic data.

VCF files store the variation data detected in one or more individuals, ensuring that the file sizes are reduced, but detailing the specific alleles that are different from the reference sequence.

All file types can be indexed creating a separate index file, which is a requirement for distributing the files by htset.

### 4.3.6 Crypt4GH

Crypt4GH<sup>11</sup> is an encryption standard that uses one or more encryption keys to encrypt one or more regions of a file, such as a VCF or SAM / BAM / CRAM files, but can be used for any file type of format. The encryption key or keys that are used to encrypt the file itself are added to the header of the file encrypted with secondary keys. To allow decryption of one or more regions of the file, the user just needs the specific key or keys to decrypt the portion of the header which encrypted the region of interest in the file, and then use that key to decrypt the required region. This has the advantage that encrypting the file for different users just needs the specific keys to be encrypted per user, reducing re-encryption overhead. It also allows only specific regions of a file to be decrypted for certain users, restricting access to certain regions or genes.

### 4.3.7 Authentication and Authorisation Infrastructure and Passports.

The required Authentication and Authorisation Infrastructure (AAI) and Passport standards are defined by the GA4GH. These standards specify the way a user's identity and access permissions can be distributed and understood across a federated network. The standards also allow a user to link multiple identities together. For example, the ELIXIR AAI<sup>12</sup> conforms to the GA4GH AAI<sup>13</sup> standard, and supports the GA4GH passport<sup>14</sup> standard. This allows users who have an ELIXIR identity to link other identities to their ELIXIR identity. This allows individual nodes to maintain their own identity provider, if required, but still allow a user to use their ELIXIR identity to access authorised resources across a federated network. For example, a user can link their ELIXIR identity to an EGA identity, which would allow the user to access the Resource Entitlement Management System<sup>15</sup> (REMS) hosted by CSC<sup>16</sup> to administer user permissions using their ELIXIR identity, and then access the EGA using the

---

<sup>11</sup> <https://doi.org/10.1093/bioinformatics/btab087>

<sup>12</sup> <https://elixir-europe.org/services/compute/aa1>

<sup>13</sup> <https://github.com/ga4gh/data-security/tree/master/AA1>

<sup>14</sup> [https://github.com/ga4gh-duri/ga4gh-duri.github.io/blob/master/researcher\\_ids/ga4gh\\_passport\\_v1.md](https://github.com/ga4gh-duri/ga4gh-duri.github.io/blob/master/researcher_ids/ga4gh_passport_v1.md)

<sup>15</sup> <https://www.csc.fi/en/rem-s-kayttovaltuuksien-hallintajarjestelma>

<sup>16</sup> <https://www.csc.fi/en/home>



same identity while internally the EGA identifies the user using the EGA account. The granularity of data access can be defined by the DAC depending on the ELSI requirements, for example data access can be granted to specific files, files associated with an individual, or to a cohort.

#### 4.3.8 Workflow and Task Execution Service

The GA4GH Workflow Execution Service<sup>17</sup> (WES) standard specifies how complex analytic pipelines can be specified in a standardised way to allow these pipelines to run in different locations, for example at different federated node: an important requirement in allowing users to send the analysis to the data rather than bring the data to the analysis which helps to ensure the data remains secure and conforms to the ELSI issues that apply to it.

The Task Execution Service<sup>18</sup> (TES) allows analysis pipelines to run on a wide range of heterogeneous infrastructure. For example, a user may wish to run the same variant calling pipeline on data hosted at three different nodes, but the compute infrastructure used by these nodes is different. While WES allows the pipeline to be specified in a standardized way, TES ensures that the task required to run the pipeline can be executed on different types of compute infrastructure, such as different high performance computing (HPC) or Cloud environments. WES and TES together ensure that identical, repeatable, and consistent analysis can be run on different federated node irrespective of the underlying compute infrastructure.

## 4.4 Infrastructure Technologies

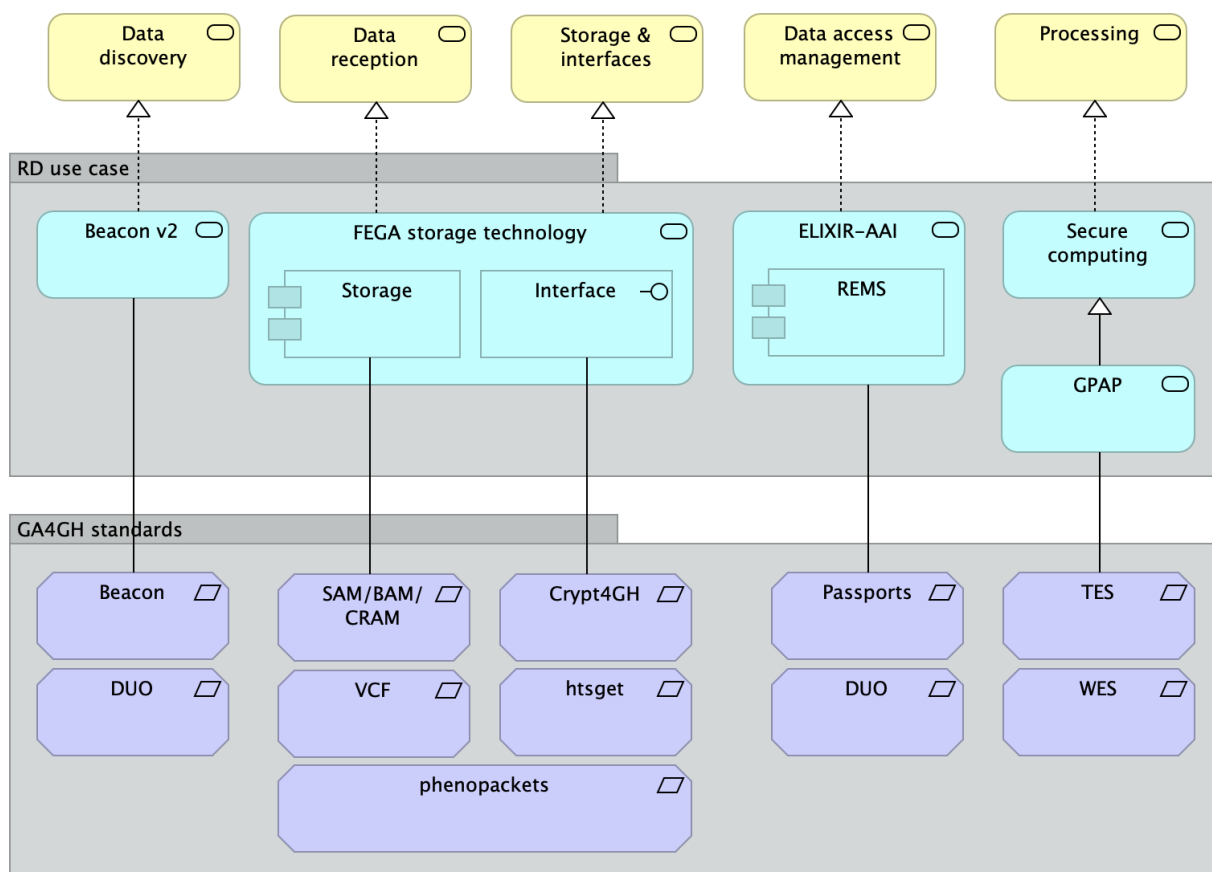
The infrastructure proposed is shown in Figure 2, which illustrates how the services and their constituent components support the five main functionalities and are supported by the range of GA4GH standards, which supports the requirements of the rare disease use case which are detailed within Section 4.6.

---

<sup>17</sup> <https://github.com/ga4gh/workflow-execution-service-schemas>

<sup>18</sup> <https://github.com/ga4gh/task-execution-schemas>





**Figure 2:** The core functionalities as defined in the scoping paper are supported by the services and GA4GH standards. The Beacon, Federate EGA technology, ELIXIR-AAI, and GPAP are all used within the rare disease use case.

## 4.5 Synthetic Data

To ensure that the proof of concept does not risk exposing real participants' data, it will use synthetic data initially to demonstrate and test the functionality of the system.

There are different requirements for synthetic datasets to meet, for example initial functionality testing, stress and capacity testing, and security and data protection testing. For the secure access demonstrator, the requirements are to demonstrate the functionalities and how they support the use cases, elicit feedback from interested stakeholders, and allow WP4 to extend and develop the roadmap over the final year of the project. The proof of concept will use a modified version of a synthetic dataset<sup>19</sup> developed for the H2020 Common Infrastructure for National Cohorts in Europe, Canada, and Africa<sup>20</sup> (CINECA) project. This dataset includes genetic data from the 1000 Genomes project<sup>21</sup>, with

<sup>19</sup> <https://ega-archive.org/datasets/EGAD00001006673>

<sup>20</sup> <https://www.cineca-project.eu/>

<sup>21</sup> <https://www.internationalgenome.org/>

phenotypic data generated from the UKBiobank<sup>22</sup> data dictionary<sup>23</sup>. The phenotypic data includes demographic data, diabetes and cardiovascular data, cancer data, and Covid data, making this dataset also applicable for the common complex disease, cancer, and Covid-19 use cases. The dataset currently does not have rare disease Human Phenotype Ontology<sup>24</sup> (HPO) or Orphanet Rare Disease Ontology<sup>25</sup> (ORDO) terms associated with it, so these terms will need to be generated. Additionally, the genetic data of the individuals will be edited to include plausible causative mutations for the associated phenotypes and disease. This will allow any analysis performed to return scientifically relevant results, ensuring that any proposed analysis is fit for purpose.

As the infrastructure develops, stress and scalability testing will be required, and to support this a larger synthetic dataset will be required which will be for loading at each data hub depending on the expected size of the data hosted by each hub. A synthetic dataset with 1 million genomes has been funded by the Finnish Ministry of Social Affairs and Health which can be utilised for stress and capacity testing of the network in 2022.

## 4.6 Rare Disease Use Case

During the B1MG Synthetic Genomes workshop, three primary rare disease scenarios were identified from the WG8 - Rare Diseases working group as detailed in Table 2.

Scenario	Context	Questions
1	An undiagnosed child with mental retardation and facial dysmorphism with a negative array-CGH analysis underwent WES, which disclosed a <i>de novo</i> , likely pathogenic, rare variant in <i>FBX011</i> gene.	<ol style="list-style-type: none"> <li>1. Is there a relationship between the gene mutation and disease?</li> <li>2. Are there any other individuals with the same mutation or allelic variant?</li> <li>3. If there are, what is their phenotype? What can be derived from them for the prognosis?</li> <li>4. What is the variant frequency across different populations?</li> </ol>
2	A girl was diagnosed affected by a Noonan-like syndrome, but she was negative after analysis of an extended panel for Rasopathies. Following WES, a <i>de novo</i> likely pathogenic mutation was found in <i>MAPK1</i> gene.	<ol style="list-style-type: none"> <li>1. Are there any other individuals with the same mutation or allelic variant?</li> <li>2. What is their phenotype and natural history of the disease?</li> <li>3. What is the variant frequency across different populations?</li> </ol>

<sup>22</sup> <https://www.ukbiobank.ac.uk/>

<sup>23</sup> <https://biobank.ctsu.ox.ac.uk/crystal/index.cgi>

<sup>24</sup> <https://hpo.jax.org/app/>

<sup>25</sup> <https://www.orpha.net/consor/cgi-bin/index.php>



3	A patient affected by spondylometaphyseal dysplasia with corner fractures, a rare form of AD osteochondrodysplasia of unknown genetic origin underwent WES analysis, which disclosed a <i>de novo</i> likely pathogenic mutation in <i>FN1</i> , a gene previously associated with glomerulopathy with fibronectin deposits, a rare kidney disease (KD) not present in the patient.	<ol style="list-style-type: none"> <li>1. Are there any other individuals with the same mutation or allelic variant?</li> <li>2. What is their phenotype?</li> <li>3. Is it possible to establish any genotype-phenotype correlation?</li> <li>4. What is the variant frequency across different populations?</li> </ol>
---	---	--

**Table 2:** 3 example use cases as defined by WG8 and given to WP4 as example use cases.

By removing duplicated questions from the 3 use cases above, we can determine a minimal set of questions that would support all 3 use cases, and map these to the required user-facing services (Table 3).

Question	Primary Functionality	User facing Service
Search for individuals with a particular variant.	Data Discoverability	Beacon
Determine the frequency of specific genes with variants in a particular node or set of nodes.	Data Discoverability	Beacon
Explore variants from selected individuals.	Processing	Genome-phenome Analysis Platform
Explore phenotypes / disease of selected individuals.	Processing	Genome-phenome Analysis Platform
Search for individuals with variant(s) in a particular gene or region.	Data Discoverability	Beacon V2
Search for individuals with a particular phenotype or disease.	Data Discoverability	Beacon V2
Determine the frequency of variants, or specified genes with variants, in a particular population.	Data Discoverability	Beacon V2
Determine the frequency of variants in a particular node or set of nodes.	Data Discoverability	Beacon V2
Return the identifiers of the individuals with a certain phenotype and variant, or certain phenotype and variants in a specified gene, or a certain disease and variant, or a certain disease and variants in a specified gene.	Data Discoverability	Beacon V2
Return a set of variants called on different individuals in different nodes which were called using the same analysis pipeline to ensure data harmonisation.	Processing	WES, TES

**Table 3:** List of questions and associated functionality, service and type of goal in the PoC



As described above in the GA4GH section, the Beacon V2 has not been submitted to GA4GH for approval, and hence the approval of this new standard is required to ensure the full functionality of the proof of concept. Hence a distinction has been made between the existing Beacon standard, and the rare disease use case questions that it can answer, versus the questions which require the extended functionality of the Beacon V2, for which a stretch goal has been defined.

## 4.7 Other use cases

Over the next 12 months, it is proposed that WP4 elicits feedback from the other working groups to extend the requirements for the core infrastructure to the cancer, infectious disease, and common complex disease use cases. By working with the CINECA project, specifically WP4 and WP5 (federated analysis and cancer use case respectively), B1MG hopes to leverage their developments for B1MG. For example, the inclusion of the federated analysis stretch goal for the rare disease use case is in response to the CINECA WP4 federated variant calling pipeline which utilises the GA4GH TES standard. B1MG is working closely with CINECA for their GWAS use case, investigating the federated genotype / phenotype association analysis using GWAS.

## 4.8 Data security and GDPR

To maintain trust in a system that allows federated analysis of genomic and phenotypic data, it is crucial that the data is maintained in a way that respects the nine core ELSI principles as defined by WP2. To this end, an initial gap analysis was done with Ethical, Legal, and Societal Issues (ELSI) and associated requirements, and the main General Data Protection Regulation<sup>26</sup> (GDPR) principles as supplied by WP2 to WP4, and the functionalities, components, and services proposed for the core infrastructure and to be implemented as part of the PoC. The initial results are summarised below.

A bottom up approach was initially taken to try and identify any services or associated procedures or processes (where processes are a constituent process of a large component or service) that may not conform to data protection principles (DPPs) or GDPR as part of the Data Protection by Design and Default approach described by WP2<sup>27</sup>. This enabled the creation of a heat map (Figure 3), which detailed the data protection principles which were fully met, needed clarification or supporting documentation or protocols, and not met. For example with the logging process within the PoC, this conforms to the accuracy DPP but

<sup>26</sup> <https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:32016R0679>


<sup>27</sup> <https://drive.google.com/file/d/1nhYP2aSX9v3dAqxwdOM0olegUtFz2DTv/view?usp=sharing>





procedures regarding the length of time these logs should be kept for (Storage Limitation) need to be clarified when these logs contain data subject to GDPR.

Service, process, or procedure	ELSI Principles								
	Transparency	Lawfulness	Fairness	Purpose Limitation	Data Minimisation	Accuracy	Storage Limitation	Integrity and Confidentiality	Accountability
Logging	Yellow	Yellow	Yellow	Yellow	Yellow	Green	Yellow	Yellow	In progress
Metadata	Yellow	Yellow	Yellow	Yellow	Green	Green	Yellow	Yellow	
Archive	Yellow	Yellow	Yellow	Yellow	Green	Yellow	Yellow	Yellow	
Distribution	Yellow	Yellow	Green	Yellow	Green	Green	Green	Green	
Submission	Yellow	Red	Green	Yellow	Red	Red	Yellow	Yellow	
Access Restriction	Yellow	Green	Green	Yellow	Green	Green	Yellow	Green	
Discovery	Yellow	Yellow	Yellow	Yellow	Green	Green	Yellow	Yellow	
Data Linkage	Red	Red	Green	Green	Green	Green	Green	Green	



- Service complies, or is based on another service
- Service requires legal definition or agreement of terms, or clarification to comply
- Service does not currently comply

**Figure 3:** Heatmap of how core services, processes and procedures comply with ELSI principles.

Initial DPPs selected for prioritisation during the workshop include data minimization and purpose limitation. Data minimization means here that only data needed for a certain project will be selected and accessed. This applies to both datasets of data subjects as well as data types or parts of the genome relevant for the envisaged use. On the practical level, this is related to purpose limitation as this DPP also requires a limitation of accessed data, here determined by the data use conditions as defined. Both, again, require as a safeguard also access restriction to as few people as possible. In addition, safeguards to ensure purpose limitation need to make sure that neither malicious nor accidental processing for other purposes can take place, thus having a focus on the security setup of the processing. These two DPPs will be focussed on initially to ensure that all services, procedures, and documentation are identified to ensure these are fully met by the PoC.

WP4 will work to address the services or components that clearly fail to meet these DPPs, and work with WP2 to gain clarification where required on the requirements, and to detail the procedures or documentation that is needed to ensure compliance where required. In addition to the bottom up approach, we will use the PoC to enable a top down approach, using the PoC to determine how the PoC meets the DPPs from a user perspective, and use this to test specific use cases, such as a user wanting to withdraw consent and needing information on how to initiate the process, or a submitter wanting to ensure the data hub meets their responsibilities as a data controller.

WP4 will also reach out to other projects, such as CINECA, Solve-RD, EJP-RD, and Federated Human Data<sup>28</sup> to determine how these projects are meeting the DPPs, and if there are existing solutions that can be utilised within the 1+MG and B1MG projects. Extensive information is required for the nodes on how to set-up and run a data hub that conforms to all the DPPs, so we will strengthen our collaboration with Federated Human Data to identify, document, and share these requirements and procedures with the prospective data hubs. This documentation will also identify where the responsibility lies for each service or component, or procedure that is needed to ensure the data hub meets the DPPs. The use of the initial size limited synthetic dataset will allow the data hubs to test the processes while not breaking any of the DPPs.

Along with basing the infrastructure and components on GA4GH standards, which have been approved from a Data Security and Regulatory and Ethics respect before approval by GA4GH (Appendix 1), WP4 will use the PoC to engage with WGs 8 - 11 and WP2 to extend the gap analysis to all use cases. The PoC will also allow WP4 to define the Standard Operating Procedures (SOPs) required for each node hosting the infrastructure to ensure compliance with GDPR and to maintain the security of the hosted data. While WP4 defines the infrastructure required to support secure cross border data sharing, WP4 does not have the resources or responsibility to ensure that the hardware and physical infrastructure used to support the components and services recommended by WP4 are suitable and comply to the requirements set by GDPR, WP2, the WGs 8 - 11, and the data controllers. As the use-cases give feedback on the proposed infrastructure WP4 can define and clarify the responsibilities of the nodes with respect to security and data protection.

## 4.9 Planning

Figure 4 details the high level roadmap, presented as a Gantt chart, for the remainder of 2021, for the ELSI issues, the development of the synthetic datasets required to demonstrate the PoC, and the technical roadmap for deploying and running the infrastructure required for the PoC. There are a set of dependencies which need to be met to allow the PoC to be deployed, which determine the timelines.

---

<sup>28</sup> <https://elixir-europe.org/about-us/commissioned-services/federated-human-data>



	Q2	Q3	Q4
<b>ELSI</b>			
Complete Gap Analysis (bottom up)	█		
Document ELSI responsibilities (service / standard / node)	█	█	
Outreach to other projects	█	█	
Outreach to other WGs			█
Iteratively refine infrastructure based on GDPR requirements		█	█
<b>Synthetic Data</b>			
Spike raw data	█		
Spike aggregate data	█	█	
Add phenotype / disease associations	█	█	
Generate distinct datasets for nodes			█
<b>Technical requirements</b>			
Service level requirements documented	█		
Outreach to initial nodes	█	█	
Feedback of technical requirements		█	
<b>PoC deployment at Early Adopter Nodes</b>			
Storage and Interfaces	█	█	
Data Discoverability	█	█	█
Data Access and Management Tools		█	
Data Reception		█	█
Processing		█	█
Synthetic Data Load / Testing		█	█
<b>PoC Deployment at Second Wave Nodes</b>			
Infrastructure Deployment		█	█
Synthetic Data Load / Testing			█

**Figure 4:** The remaining high-level tasks needed to ensure the PoC is delivered on-time by Q4 2021.

Figure 5 defines the proposed high level roadmap for 2022. After demonstrating the rare disease PoC to the use cases and signatory member states, we will define any gaps in the requirements from the Covid-19, common disease, and cancer use cases not covered by the PoC. This will include extending and adapting the synthetic dataset for the remaining use cases, for example by adding sets of variants with a known association to a common disease, such as diabetes, or cancer related genes. The synthetic dataset will need to be

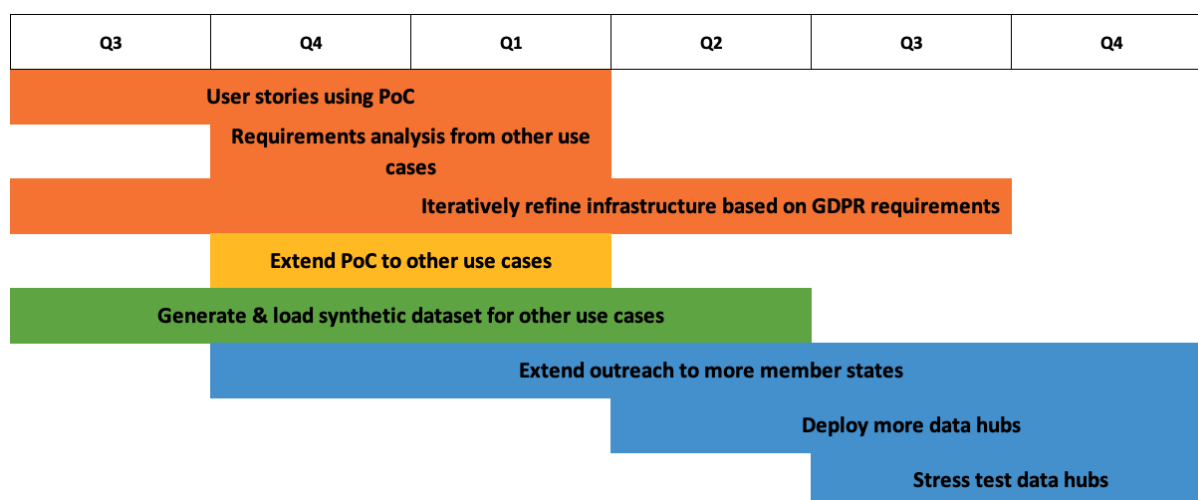


scaled up, to allow stress testing of the complete network, and split into more datasets applicable to the number of individuals each data hub is expected to host. An expansion of the outreach will be needed to engage with the remaining signatories to demonstrate the infrastructure and elicit their expected data volumes and timelines.

Included in the outreach will be the results of the ELSI work, defining the SOPs and responsibilities needed to ensure the data hub corresponds to the DPPs, utilising the information gained by installing, deploying, and running the PoC demonstrator.

Collaboration will continue with relevant projects, such as CINECA, Federated Human Data, ELIXIR-CONVERGE, and Federated EGA to maximise leverage work in these areas. For example, Federated EGA are writing the documentation and procedures required to operate a Federated EGA node; it is expected that some of the procedures will be applicable to the B1MG data hubs.

Close collaboration with GA4GH will continue to be required through the ELIXIR and GA4GH strategic partnership, to ensure the standards already deployed are continue to be developed, maintained within the B1MG infrastructure, and kept secure but also to develop the cloud compute processes required to enable more secure federated analyses to be performed on the data as required by WGs 8 – 11.



**Figure 5:** Outline of the plan to Q4 2022 - extending the PoC to other use cases and iterating the infrastructure definition based on GDPR requirements, gaining additional requirements, generating and loading the final synthetic datasets, and deploying the data hubs to the member states.

## 5. Results

The document describes the proposed evolving roadmap towards a demonstrator for secure cross border data access by M24.

## 6. Discussion

This document allows all interested stakeholders to understand the proposed timelines for the PoC demonstrator, and therefore at what point their requirements or responses will be required



to continue with the development of the roadmap. As the roadmap is a living document input is encouraged from all stakeholders to ensure the document covers as many use cases as possible.

## 7. Conclusions

The document delivers the first version of the roadmap to provide the infrastructure to support cross border data access of restricted access human genetic and phenotypic data.

## 8. Next steps

Continue to engage with stakeholders, especially WG 8 - 11, to continue to evolve the roadmap to ensure it tracks the required development to ensure the demonstrator supports all use cases.

## 9. Impact

The roadmap, along with the initial scoping paper, establishes the prerequisites and timelines required to provide the demonstrator, and the timings for use case interaction.

## 10. Appendices

### Appendix 1

Once a proposed product has been incorporated into the GA4GH roadmap, the product is submitted to the following bodies for review:

1. The Data Security Foundational Workstream (DSWS),
2. the Regulatory and Ethics Foundational Workstream (REWS),
3. and a specially convened Product Review Committee (PRS).

The PRC consists of three members who are nominated by the submitter but must comprise of one representative from each of the groups below:

1. A workstream leader from another technical workstream,
2. a member of a third, different, technical workstream,
3. and a representative of a driver project which has been involved in the product development.

The PRC must return a unanimous 'Accept' verdict for the product to pass this stage in the approval process. The PRC may also refer the product back to the DSWS or the REWS with any concerns they may have. Additionally, the REWS and DSWS must both 'Accept' the product.



Once all three bodies have made positive assessments of the product, the product is presented to the GA4GH Steering Committee, where the steering committee will vote on whether to approve the product or not. Rejection will require re+submission, though reasons for rejection will be given. Approval by the steering committee will formally make the product an approved GA4GH standard.

Once a standard has been approved, changes to the standard must be versioned using semantic versioning. This ensures that interoperability is maintained within the standard while allowing the standard top to evolve as technology and requirements change.

By utilising GA4GH standards where possible, the B1MG infrastructure ensures that the standards have gone through an approval process which addresses security, ethical, and legal issues. Additionally, the approval process ensures that the standard does not duplicate, or replicate the functionality of an existing standard, reducing the proliferation of standards. Major changes require that the standard is re-submitted to GA4GH for approval.

