

# Towards federated service discovery and identity management in collaborative data and compute cloud infrastructures

Shiraz Memon, Jens Jensen, Willem Elbers,  
Helmut Neukirchen, Matthias Book and Morris Riedel

## Published version information

**Citation:** S Memon et al. "Towards federated service discovery and identity management in collaborative data and compute cloud infrastructures." Journal of Grid Computing, vol. 16, no. 4 (2018): 663-681.

**DOI:** [10.1007/s10723-018-9445-3](https://doi.org/10.1007/s10723-018-9445-3)

*This is a post-peer-review, pre-copy-edit version of an article published in Journal of Grid Computing. The final authenticated version is available online at DOI above.*

This version is made available in accordance with publisher policies. Please cite only the published version using the reference above. This is the citation assigned by the publisher at the time of issuing the AAM. Please check the publisher's website for any updates.

# Towards Federated Service Discovery and Identity Management in Collaborative Data and Compute Cloud Infrastructures

Shiraz Memon  · Jens Jensen  · Willem Elbers  ·  
Helmut Neukirchen  · Matthias Book  · Morris Riedel 

Received: date / Accepted: date

**Abstract** This paper compares three multi-national research infrastructures, one that provides data services, one that provides compute services, and one that supports linguistics research. The aim is to jointly provide services to the user communities, and, perhaps eventually, seamlessly interoperate. To this end, we look at and compare how the infrastructures build their service federations (trust, service status, information systems), and how they manage users (identities, authentication, and authorisation).

**Keywords** Distributed infrastructure · Federated identity management · Service discovery · Standards · Interoperation · Cloud computing

## 1 Introduction

Distributed compute, data, and more recently, cloud infrastructures have been successful in providing resources to a wide variety of research communities. The e-Infrastructure Reflection Group identified in 2004 the

outline/vision of a distributed infrastructure comprised of fabric (disk, CPU, networks), and a “middleware” layer connecting the infrastructure across sites; user communities would then develop and deploy their own applications on top of the e-infrastructure [44]. Also the Foster/Kesselman vision of grid computing [31], with computing available on demand through standard interfaces, was hugely influential in the development and use of e-infrastructures, leading for example to the middleware that is known as Globus Toolkit [29] and more recent Globus cloud services [30].

The established e-infrastructures have been very successful, having provided resources to researchers on a national or multinational scale in TeraGrid [36], European National Grid Initiatives (NGIs), Extreme Science and Engineering Discovery Environments (XSEDEs) [52], or, in the case of the world-wide Large Hadron Colliders (LHCs) Computing Grid, a truly global scale [45]. They have provided data and compute resources in support of a vast range of research.

The main contribution of this paper is *connecting the infrastructure*, particularly focusing on security and service discovery (Fig. 1). There is plenty of existing work on e-infrastructure architecture and security, managing users and their communities [2,8,13,18], which we summarise below for the reader’s convenience. We are, however, interested in the practical applications, so we have chosen three infrastructures with different purposes and look at the general challenges of bridging them, as well as connecting their user communities. We also look at the specifics of some of the key services involved in this endeavour, going into details of recent developments.

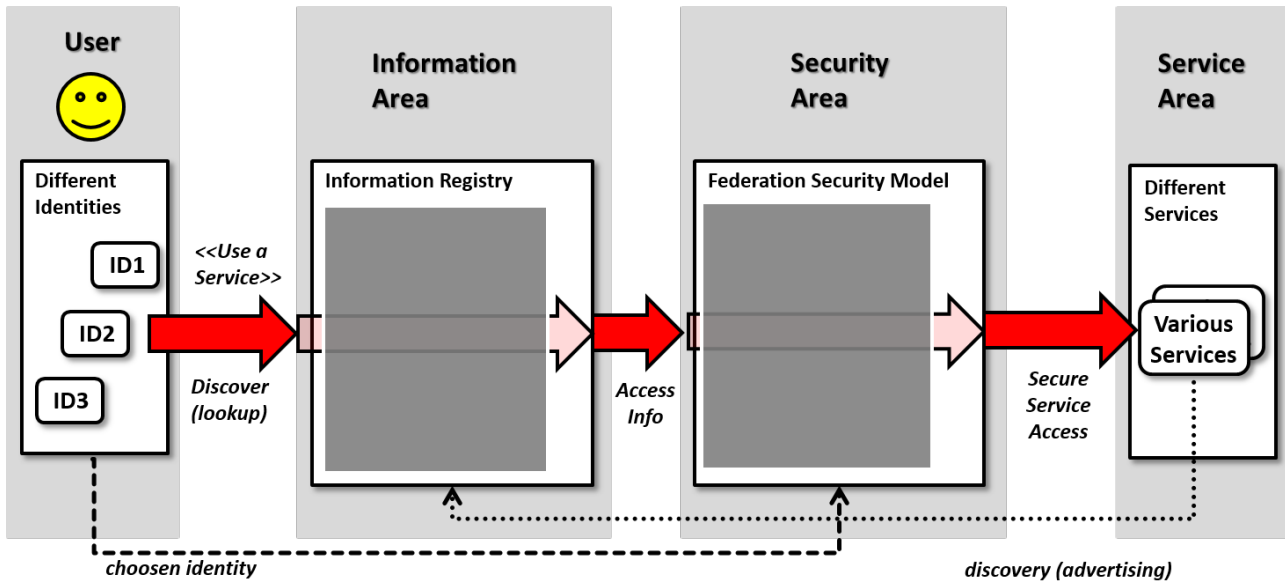
---

S. Memon · H. Neukirchen · M. Book · M. Riedel  
University of Iceland  
Reykjavik, Iceland  
E-mail: helmut@hi.is, book@hi.is

S. Memon (✉) · M. Riedel  
Jülich Supercomputing Centre, Forschungszentrum Jülich  
Leo-Brandt Straße, 52428 Jülich, Germany  
E-mail: a.memon@fz-juelich.de, m.riedel@fz-juelich.de

J. Jensen  
STFC, Harwell Oxford Campus  
Didcot, United Kingdom  
E-mail: jens.jensen@stfc.ac.uk

W. Elbers  
CLARIN ERIC  
Utrecht, Netherlands  
E-mail: willem@clarin.eu



**Fig. 1** EUDAT federated authentication and service discovery. The EUDAT architecture without any specific authentication and service discovery architecture. On the right are the EUDAT's B2\*services.

### 1.1 Connecting the Infrastructure to Itself

The following components are the key components to defining and binding together an infrastructure:

- Common fabric security, i.e., X.509 host certificates from trusted Certification Authorities.
- Service naming: Each relevant service must have a *name* by which it can be discovered and referenced; a typical type of name is a Web services endpoint or URI.
- Service discovery/metadata: a way to discover which services would be available to the user.
- Service registry: a location where each service is registered, typically used to record whether it is a legitimate part of the infrastructure and whether there are scheduled downtimes, etc.
- Service information granularity: The information model representing the service should be sufficiently flexible to capture the service details from a coarse-to fine-grained level. Furthermore, the model must be interoperable as multiple infrastructures are discovering and advertising their services.
- Operations and support: From the user's perspective, there should be a single point of contact for support, and there should be a *team* responsible for operating the service (as opposed to individual admins at each site.)

### 1.2 Connecting Users to the Infrastructure

Central to the e-infrastructures that are a focus in this paper are:

- Common authentication: This allows each user to access any part of the infrastructure with a single credential (as well as accessing other infrastructures with the same credential);
- Service discovery mechanisms: There has to be an “entry point” which helps users discover services that are available to them. Typically, this is a portal, but could also be hosted on a “user interface” node (to which users log in or connect with remote desktop);
- Service database (which may or may not be the same as the service discovery): Typically, it is a central database listing the services that are part of the e-infrastructure. By extension, an associated service could be used to monitor service status, announce scheduled downtimes, etc.;
- Common authorisation: This is needed across the infrastructure to provide additional actions to researchers and users, enabling them to share data and to collaboratively make use of the services provided.

Note the difference between the service discovery and registry/database in Sect. 1.1 and 1.2: While they might be the same service in some infrastructures, the former is more likely to have an Application Programming Interface (API) to allow programmatic access (cf. R14 below), or technical interfaces for administrators, whereas

the latter should be browser-accessible and more user-friendly.

## 2 Architecture and Concept Backgrounds

Unsurprisingly, the e-infrastructures covered here are architecturally similar; even with independently designed architectures they end up often providing the same types of services. Indeed, one of the achievements of the AARC project was a unified view of the authentication and authorisation parts of the e-infrastructures [32]. Also, common standards and interoperation play an important role, such as the GLUE standard (Sect. 4.3.1), as they enable service discovery across domains if used correctly [19].

Table 1 shows an overview of how the three different infrastructures provide interfaces for their users and how they are connected internally. Here, “CLI” is short for “command line interface” (which is generally considered harder to use for novices but saves time for experts); “WS” refers to web services for programmatic access; and X.509 is the standard for certificates [18] provided through IGTF ([www.igtf.net](http://www.igtf.net)). VOMS is the Virtual Organization Membership Service, an attribute authority [2]. Finally, BDII (Berkeley Database Information Index) and GOCDB (Grid Operations Centre DataBase) are information services, used for service discovery and registry, respectively, and are covered in more detail in Sect. 3.3.

## 3 Requirements Analysis

In today’s research environments, Single Sign-On (SSO) is an important requirement: It enables researchers to use a single account to access remote services, and service providers do not need to maintain separate account data, nor do they need password quality checking, password reset, maintaining user contact details, etc. Importantly, researchers present the same identity and can use the same credential with several different services, so SSO can potentially bridge infrastructures.

Extending SSO, national research networks build identity management federations where Identity Providers (IdPs) are bound by common federation policies, thus ensuring a common level of assurance (LoA) of identities and a common set of attributes being passed to the services. These attributes are used to identify (or at least represent) the user to the service, and/or used for authorisation. Typically, these national federations use web-based technologies (users use a

web browser to access services via portals), such as the SAML Web Single Sign-On (Web SSO) profile, and use (subsets of) the eduPerson schema to publish attributes.

As much research is international, it becomes useful to connect national identity federations, despite their publishing different attributes or having different levels of assurance (LoAs). eduGain [20] is an inter-federation identity management framework, which aims at interconnecting the national federations. However, there is still a need for harmonisation due to the differences between national federations; this is the subject of ongoing work from REFEDS ([www.refeds.org](http://www.refeds.org)) and recent work from the Authentication and Authorisation for Research and Collaboration (AARC) project [1]. As we shall see, one option for infrastructure projects is to implement a *proxy* to harmonise credentials [14], and perhaps, via *credential translation*, provide support for non-web (command line) access. The other main option is to simply implement a project or community-specific independent (non-federated) IdP. Obviously, many of the advantages of SSO are then lost, but as we shall see, the adherence to standards creates opportunities for interoperation between infrastructures.

In the following subsections, we analyze the requirements from three different infrastructures: a research community infrastructure, a data infrastructure, and a compute/cloud infrastructure, the latter two being multi-disciplinary. We look at these as individual infrastructures (cf. Table 1), but also at how they can share users and services such as workflows.

### 3.1 CLARIN European Research Infrastructure

Common Language Resources and Technology Infrastructure (CLARIN) [15] provides easy and sustainable access for scholars in the humanities and social sciences to digital language data (in written, spoken, or multi-modal form), as well as access to advanced tools to discover, explore, exploit, annotate, analyse or combine the data, regardless of where it is located. CLARIN is building a networked federation of language data repositories, service centres and knowledge centres, with SSO access for all members of the academic community in all participating countries. Tools and data from different centres are interoperable, so that data collections can be combined and tools from different sources can be daisy-chained to perform complex operations.

The CLARIN infrastructure is fully operational in many countries, and a large number of participating centres are offering access services to data, tools and expertise. At the same time, new services are added by countries that joined more recently, and CLARIN’s

	Service	CLARIN	EUDAT	EGI
User	Authentication	federated/own†	federated/own†	X.509/federated/own†
	Access methods§ (Web/CLI/WS)	Web	Web/CLI	CLI/WS
	Authorisation	own	own	VOMS
	Service discovery	Portal/Switchboard	Portal	Wiki
Infra	Workflow	WebLicht	N/A	N/A
	Authentication	IGTF	IGTF	IGTF
	Service discovery	Portal/Switchboard	N/A	BDII
	Service registry	Switchboard	GOCDB	GOCDB

**Table 1** Infrastructures need ways to give access to users, and to link services within the infrastructure. Some are the infrastructure’s own, others are shared or come from an external federation. Abbreviations are explained in section 2.

datasets and services are constantly updated and improved. On the services page [16] we show the services accessible at this moment, and explain how and by whom the various services can be accessed.

### 3.1.1 Requirements

- R1 *Single Sign-On (SSO)*: To provide single sign on, users must be able to use a single identity for all CLARIN services, and credentials should only be required for the first authentication. Authorization within the CLARIN infrastructure is not centrally managed, but on a service per service basis. This is a result of the distributed nature of the infrastructure, where each CLARIN centre is responsible for the services it runs.
- R2 *Delegation* of user rights is crucial in a distributed service oriented infrastructure such as the CLARIN infrastructure [11]: A user typically stores data in a workspace and wants services, possibly hosted at other centres, to process the data in these workspaces. The user is authenticated and authorized to the service and then wants to delegate his/her identity and permission to the service, so the service can access the workspace on behalf of the user.
- R3 *Service discovery*: Given a dataset, what services are available to process this dataset? Given a service, what other services are available to operate on the output of this service? It is necessary to have a discovery service which describes the services’ capabilities and provides endpoints for accessible resources and services. An example of such a registry is the Language Resource Switchboard [53]. It is important to point out that such a service registry is not a workflow composition engine itself; instead a workflow composition engine typically queries a service registry during workflow composition.

### 3.2 EUDAT

European Data Infrastructure (EUDAT) [26] is a European data infrastructure which facilitates management and federation of “big (research) data” across Europe. It operates a number of services to deposit, replicate, and archive data. Services are geographically distributed across different organisations (which are currently the same as the project partners).

#### 3.2.1 Requirements

- R4 *Single Sign-On (SSO)*: Users should be able to access EUDAT services while authenticating with their “home” credentials issued by their organisation’s identity provider. Without SSO, the users would have to register with every service, and each service in EUDAT would have to maintain its own user database. This would not be scalable and might lead to inconsistencies where the same information is stored in multiple databases. Therefore, the Authentication and Authorisation Infrastructure (AAI) technology must be able to support SSO.
- R5 *Distributed authorisation*: Once users can *authenticate*, the infrastructure needs to provide an *architecturally central* authorisation service (i.e., there is only one) which is consistently enforced across the *distributed* services. The main goals are: (1) harmonised authorization policy management per service, (2) authorization decisions must be applied even in case any centralised service is unavailable, and (3) based on standards such as eXtensible Access Control Markup Language (XACML) [42].
- R6 *Non web-based federated access*: While web-based services are used as “high-level” access points, there is sometimes a need to support command line tools and “delegated” credentials. Typically these drive services based on the data transfer protocol GridFTP [3], the storage service based on iRODS [17], or services offering REST APIs.

- R7 *Delegation* of rights to other users or services is important in a data management pipeline where the service or user should be able to perform a task on a behalf of the user/owner of the data or resource.
- R8 *Multiple authentication protocols*: None of the EUDAT services were written from scratch; they were all developed around existing software products. However, there was no single authentication mechanism supported by all these products, so EUDAT's choice was to either choose a common mechanism and implement it in all services, or alternatively support multiple authentication mechanisms within the infrastructure. EUDAT, building on previous experiences in its project phase one, chose the latter. Hence, the AAI should act as an intermediary (a *proxy* in [10]) between the user and the services and *translate* the credentials from one form to the other to enable seamless access to the service.
- R9 *Different level of assurance (LoA)*: Often, most of the users perform less sensitive operations, for example reading a data set from B2SHARE (the EUDAT data sharing service). For some of the users, a high LoA is needed to perform privileged operations, for example uploading a dataset or invoking a data archival operation. A low LoA is rather useful for the volunteer scientists (e.g., holding social identities [34]) who are only interested in, say, visualisation of data. Therefore it is highly desirable for the EUDAT AAI to support segregating the service actions into different levels, hence associating each credential with a different LoA.
- R10 *Service discovery*: EUDAT infrastructure is comprised of many distributed heterogeneous services and resources: storage resources, their providers, data services, and authentication services, etc. Thus, it is essential to know the offered capabilities, types, and other specific characteristics (e.g. data transfer rate or storage capacity) of services. The infrastructure should enable users as well as monitoring system to discover the services based on the service properties.

### 3.3 EGI

The European Grid Infrastructure (EGI) [21] is one of the largest multidisciplinary grid and cloud infrastructures in Europe, hence a wide number of scientific user communities and resource providers are involved. EGI offers a set of distributed services which enable users to execute complex computing workflows. The authentication and authorisation infrastructure is based on Public Key Infrastructure (PKI), the service discovery

is supported by incorporating Berkley Database Information Index (BDII) [9] and Grid Operations Centre Database (GOCDB) [37]. This section focuses on federated authentication [14] and service discovery [22] requirements of the EGI infrastructure.

#### 3.3.1 Requirements

- R11 *Single Sign-On (SSO)*: The users should be able to use their single institutional identity to access the EGI services. Since EGI is based on PKI, users normally authenticate with their end-entity X.509 certificate. SSO will require a proxy generating a temporary certificate on behalf of the user (via a trusted online Certification Authority).
- R12 *Non web-based federated access*: Most of the EGI services are accessed through web portals, but some of them offer command line access.
- R13 *Delegation*: Users often submit compute jobs or workflows to the EGI High-Performance Computing (HPC) or cloud resources, and the user job may need to stage-in or stage-out data to a storage resource. Consequently, the compute service may need delegated access to the storage resource on the users' behalf. The delegation of rights is essential in the given use case, and in some cases credential translation is necessary as the services may not necessarily use the same authentication protocol.
- R14 *Service discovery*: In addition to providing lists of services for users and administrators, the service registry plays a significant role in composing as well as executing workflows.

### 3.4 Specific Service Discovery Requirements

- R15 *Common service information model*: The federated infrastructure registry should be able to provide a means of publishing information in a standard- and middleware-agnostic manner.
- R16 *Unified service registration and query protocol*: EGI uses different middlewares (UNICORE, ARC, Globus, HTCondor, etc.), each potentially with its own native information system. While a provider only needs to talk to their "local" information system (R18), it would be nice if all information systems had a consistent API.
- R17 *Service lifecycle management*: It is necessary to have a consistent API to manage the whole lifecycle of the services by the service providers/publishers – registration, discovery, query, downtimes, suspension, deregistration.
- R18 *Support for registry hierarchies*: Each domain (NGI) has its own registry since it can act as an infrastruc-

ture in its own right. Support for a registry hierarchy provides a unified registry for the infrastructure.

**R19 Replication of service information:** To achieve robustness within the service discovery infrastructure, the technology should support replication of information across distributed entities whereby the failure of one registry node should not hamper the functioning of other registry nodes. Moreover, better performance can also be achieved by routing traffic to less occupied registry nodes. The registry should be able to replicate its state across other registry nodes in an automated fashion.

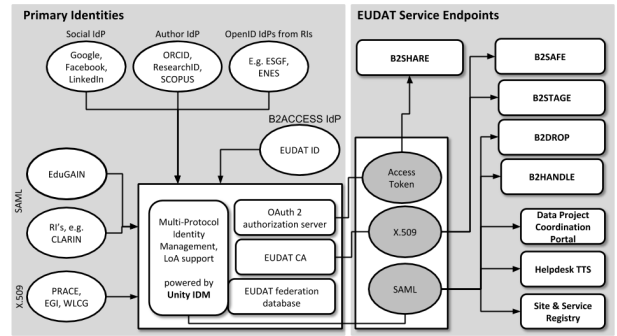
**R20 Scalability:** The registry should be able to cope with the discovery of large numbers of services in a global scale infrastructure. Since the number of services can also grow dramatically, the underlying database technology should be capable of distributing the service records horizontally and in a cost-effective manner.

### 3.5 Discussion

	Requirements	CLARIN	EUDAT	EGI
AAI	SSO	✓	✓	✓
	Delegation	✓	✓	✓
	Non-web federated access	✗	✓	✓
	Multiple authentication protocols	✗	✓	✓
	LoAs	✗	✓	✓
	Distributed authorisation	✗	✓	✓
SD	Service discovery	✓	✓	✓
	Unified API	✓	✓	✓
	Replication	✓	✓	✓
	Hierarchies	✗	✓	✓
	Service info. lifecycle management	✓	✓	✓
	Common information model	✓	✓	✓

**Table 2** Summary of the requirements analysis

Table 2 summarises the AAI and service discovery (SD) requirements from EUDAT, CLARIN, and EGI. It can be observed that most of the requirements are overlapping with each other. This, however, creates a strong motivation for having a common framework for federated service access and discovery. Although they seem similar, it is pertinent to consider certain factors, such as the number of users and services, types of services, cross organisational/domain/country service access, at-



**Fig. 2** B2ACCESS: EUDAT AAI federated user authentication and management components and the target B2\* services

tribute naming, data access policies, user and service provisioning, and attribute mapping.

## 4 Unified Federated Discovery and Identity Management

### 4.1 B2ACCESS: The EUDAT AAI Proxy

The B2ACCESS architecture is shown in Fig. 2. On the left-hand side of the diagram, B2ACCESS maps primary user identities, including a (sub)set of associated attributes, from external domains onto the EUDAT domain. The external IdPs can be connected to the B2ACCESS service by using different technologies: Security Assertion Markup Language (SAML), X.509 certificates, and OpenID Connect. For users without access to a suitable IdP, B2ACCESS itself can act as an IdP via a B2ACCESS-specific username and password. On the right-hand side of Fig. 2, the harmonised credential connects to EUDAT services also using different technologies, depending on the target service: SAML, OAuth2, or short-lived X.509 certificates. In all cases, credentials are managed by B2ACCESS and can be delegated to the target service (for credentials that support delegation), and need not be managed by the user at all: Only users who need command line tools need to download and manage credentials (in our case, the X.509 certificate).

In particular, B2ACCESS releases a unified set of attributes (Table 3) to the Service Providers (SPs) in the EUDAT infrastructure. The SPs can define authorization policies to grant certain permissions to a user based on the values of attributes associated with the user's identity. This is known as Attribute-Based Access Control (ABAC), as opposed to the more traditional Role-Based Access Control (RBAC). Examples are group membership, community membership, and LoAs.

B2ACCESS also provides account management, both for the users themselves and administrators. While many of the attribute values are gathered from the external IdP or during the registration process and are fairly stable, group membership can change more often and thus needs a management workflow, as well as delegated permissions (to community/group managers) in B2ACCESS. This is discussed in Sect. 4.1.3.

The B2ACCESS approach requires a one-time registration step for new users. The first time a user logs in by using B2ACCESS, the user is presented with a registration form. This allows us to require acceptance of license agreements and terms of use and, if needed, to request additional attributes. After completing this registration step, the actual mapping from the external identity onto the EUDAT identity is persisted in the B2ACCESS database.

#### 4.1.1 Example: Accessing B2SHARE and B2SAFE

As an illustration of the process described above, we look at the data sharing service B2SHARE. When authenticating to B2SHARE, the user is directed to B2ACCESS and authenticates via an IdP, say, a SAML IdP. B2SHARE supports OpenID Connect, so B2ACCESS converts the credential into a token which is presented to B2SHARE as an (anonymised) proof of identity. When it needs further attributes, B2SHARE obtains them from B2ACCESS via the “userinfo” API. We shall return to this example in Sect. 4.4.

When the user logs in, the SAML credential presented by their IdP is also converted into a short-lived X.509 credential.

B2SAFE needs an X.509 certificate. Typically, such a service is accessed through a portal, either one dedicated to the service, or as a feature in the user’s community portal. In this case, the *portal* generates the key pair and the certificate request, sends the request to B2ACCESS, and waits for B2ACCESS to return the X.509 certificate. B2ACCESS signs the certificate when the user has authenticated, and embeds relevant attributes into the certificate. For users requiring command line access (e.g. to B2SAFE), B2ACCESS can also generate the key pair and certificate itself, and let the user download both. The user then installs them locally and uses their command line tool. In its current implementation, B2ACCESS supports command line tools for services that use X.509, or for OAuth (via a bearer token).

Generally, converted credentials are only valid for a short period of time (hours instead of days), because they are managed on the user’s behalf by services, they are not held by the users themselves.

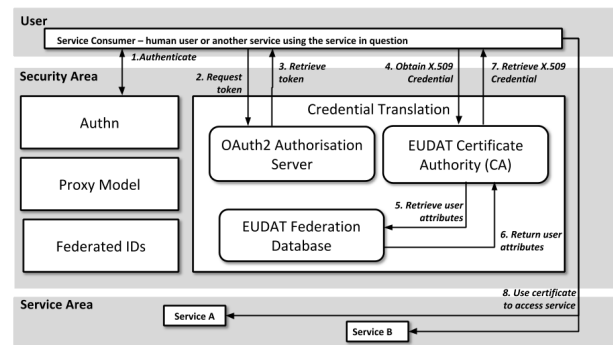


Fig. 3 Credential translation

#### 4.1.2 Attribute Harmonisation

Since B2ACCESS accepts identities from many external IdPs, and different IdPs have different attribute release policies, the incoming set of attributes is very likely heterogeneous. This makes it difficult for SPs to define authorization policies. As mentioned above, B2ACCESS acts as a proxy and tries to harmonize all incoming attribute information. This may imply mapping attribute values from other schemata onto attributes in the EUDAT attribute schema. If any essential attributes are not released by the IdP, B2ACCESS will ask the user to supply these attributes during the initial registration step.

Since users can be asked to supply values for missing attributes, and it is not considered feasible for the B2ACCESS operators to check all these values, we have concluded that a LoA per external IdP is not sufficient, but a LoA per attributes is needed, at least for the more important attributes such as e-mail or organisational affiliation. An attribute provided by a high LoA IdP gets assigned a high LoA while a user-supplied attribute value gets a low(er) level LoA. This is currently under development.

In the current implementation, however, there is a single LoA attribute, namely the LoA associated with the user’s (external) IdP (as determined by B2ACCESS operators; we do not ask IdPs to publish their LoA and would not necessarily trust the value if they did.) Typically, X.509 and Academic SAML IdPs are assigned a high level of assurance while the social and direct B2ACCESS IdPs are assigned a lower level of assurance.

#### 4.1.3 Group Management

EUDAT consists of many service providers offering a wide range of services and tools. Some of these tools are publicly accessible, but most apply authorization to at least some of the actions which can be performed



Name	Mandatory	Description
urn:oid:2.5.4.49, distinguished-Name	YES	Distinguished name (DN)
unity:persistent	YES	Persistent identifier
urn:oid:2.5.4.3, cn	YES	Common name
urn:oid:1.2.840.113549.1.9.1, userName	YES	Principal
urn:oid:2.5.4.10, o	YES	Organisational affiliation
email	YES	E-mail address
memberOf	NO	The service will perform the authorisation decision based on these roles.
loa	YES	Level of assurance

**Table 3** EUDAT Attributes

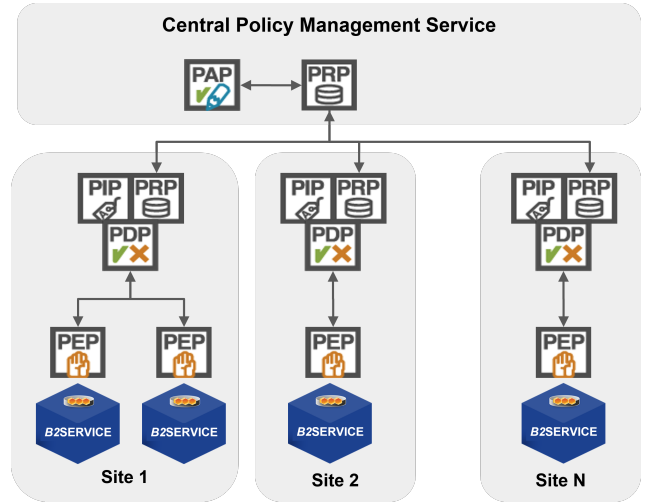
in that service or tool. As mentioned earlier, attributes released by B2ACCESS, and group membership especially, are used in these authorization policies. To provide fine-grained control, a hierarchical group structure has been defined providing: (1) a high-level domain directly under the root, defining the infrastructure, community or project, (2) multiple service level domains as children of a high-level domain, one for each service that falls under that specific high-level domain, and (3) the freedom for administrators to define anything below the service level domain to cater for any service-specific needs.

Administrators can be defined on any level to ease the administrative burden of managing the group membership. Typically, the main B2ACCESS administrators have permission in B2ACCESS to manage all groups, including the high-level domains.

#### 4.2 Distributed Authorization within EUDAT

To fulfil the requirements mentioned in R5, a solution based on XACML is under development, based on a proposed architecture shown in Fig. 4. This architecture allows for harmonised management of the XACML policies in the central service Policy Administration Point (PAP). Multiple instances of a single service can be deployed across data centres; thus, there is the need to run a central PAP and Policy Repository (PR) combination to harmonise authorisation policies for the service as a whole, covering all instances running at the individual data centres. The central service PR is replicated to the EUDAT data centres. Each data centre has a local PR. Changes are only pushed from the central service PR to the local PRs.

Each EUDAT centre is running a Policy Decision Point (PDP) with access to the local PR and each B2-



**Fig. 4** The EUDAT XACML-based distributed authorization service

service has a Policy Enforcement Point (PEP) which communicates with the centres PDP. This allows for authorization decisions even if the central service PR is unavailable.

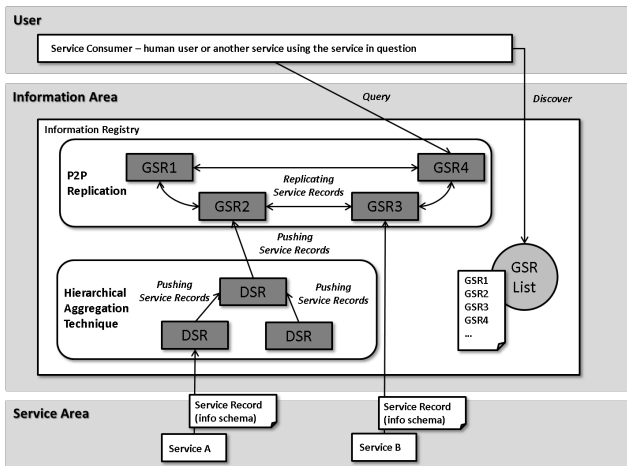
An additional function of the central PAP/PR is to provide a ingest endpoint which can be used to ingest XACML policies from external sources, such as community repositories.

#### 4.3 Federated Service Discovery with the EMI Service Registry (EMIR)

The infrastructures (EUDAT, EGI, CLARIN in our case) offer different types of services: cloud, compute, data, authentication, authorisation, etc. The EMI Service Registry (EMIR) has been designed and implemented in the European Middleware Initiative (EMI) [25] project. EMIR aims to provide robust service discovery within large scale infrastructures [27]. The initial implementation was driven by major European grid computing middlewares (UNICORE, Advanced Resource Connector (ARC), gLite, and dCache). However, the scope of the service provisioning and discovery within EMIR is not limited to grid and therefore offers a versatile service discovery utility adequate from small- to large-scale data and cloud infrastructures. The details of EMIR are described in the following subsections (see also Fig. 5).

##### 4.3.1 Concepts

The core notion of EMIR is to enable discovery of services. The set of services can be grouped in a *domain*



**Fig. 5** Federated service registry: Service discovery in heterogeneous federated infrastructure

(such as an NGI), and multiple domains can be organised in a hierarchical structure. The domain is an autonomous entity and can be connected with other domains in a hierarchy to form a *federation*. The top-level domain can replicate its information to other top-level domains in a peer-like fashion. The replication of information at the root of the hierarchy makes the federation infrastructure resilient to failures.

EMIR is based on two main components: the Domain Service Registry (DSR) and the Global Service Registry (GSR). The primary difference between the two depends on their position in the hierarchy. The DSR represents any node in the hierarchy, while the GSR always sits on the top (root node). The service to be discovered is published through a Service Record (SR) using the OGF GLUE 2.0 [47] standard.

#### 4.3.2 Service Information Model

According to the requirement R15 of a common service information model, the registry should be capable of representing the infrastructure services of any type. It could be a service having any of the (storage, cloud, network, HPC, etc.) capabilities that may dynamically (dis)appear within an e-infrastructure. In order to address the service discovery use cases from the large spectrum of scientific domains, EMIR adopts the standard GLUE 2.0 information model [47]. Since GLUE is an information model and does not provide a normative realisation, the Open Grid Forum recommendations [48] and [49] were used as a foundation to implement the service registry in the XML and JSON format, respectively. For the latter, the emerging JSON-Spec standard (similar XSD for XML) is used for the implementation. Since the GLUE model can become very extensive, in

order to be concise and yet extensible, only the abstract representation (or entities) is taken into account and forms the basis of EMIR's *Service Record (SR)*. Table 4 shows a subset of the mandatory attributes that represent a service. The JSON record in Listing 1 shows a minimal B2SHARE instance.

Attribute name	Description
Service ID	A globally unique identifier for the service
Name	Human-readable name
Endpoint URL	Location to access the service
Capability	An array of offered capabilities
Service technology	The technology used to implement the service
Service time-to-live (TTL)	The visibility of the service within an infrastructure
Service type	Service type according to namespace-based classification
Service version	Specific service version
Service health	Monitoring information about service state

**Table 4** Service record schema containing a set of core service attributes [38]

```

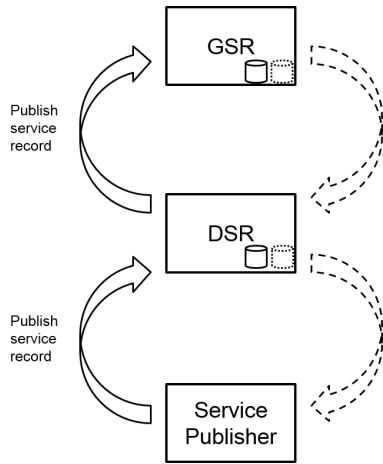
{
  "Service_Name": "B2DROP",
  "Service_Type": "eu.eudat.b2drop",
  "Service_Capability": ["data sharing"],
  "Service_Endpoint_URL": "http://b2drop.eudat.eu",
  "Service_Endpoint_Technology": "technology",
  "Service_Endpoint_InterfaceVersion": ["v1.0"],
  "Service_Endpoint_HealthState": "ok",
}

```

Listing 1: Service record in GLUE 2.0 JSON format

#### 4.3.3 Hierarchical Aggregation

EMIR allows creating flexible registry hierarchies of DSR nodes with GSR on top. Figure 6 illustrates a simplified hierarchical aggregation model where the service records are published from a leaf node (a service publisher) and traverse the DSRs to the root GSR node. The top level GSR node is *eventually consistent* [51]; however, due to the network latency of service records being published, the freshness of information could be affected. While designing the registry, two major factors must be taken into account:



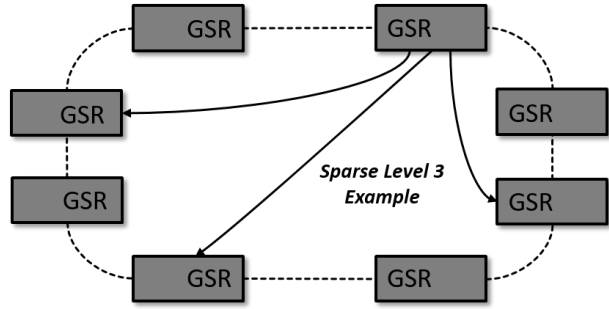
**Fig. 6** Hierarchical (bottom to top) aggregation of information

- The registries are geographically distributed across different administrative domains. In order to cope with intermediary (availability or network) failure of nodes, an in-memory database (dotted database icon in Fig. 6) is used that captures the (unsynchronised) modifications.
- A service record may contain a variety of project- or virtual organisation-specific information (apart from what has been mentioned in Table 4). Therefore, unlike conventional SQL, a schema-free or NoSQL approach (using MongoDB [39]) has been implemented. The database also offers horizontal scalability to distribute the large number of service records over multiple database instances.

#### 4.3.4 A Peer-to-Peer Approach to the Replication of GSRs

The notion of replication of GSR top level registry nodes in a hierarchy is based on the Pastry algorithm [46] and inspired by the ISIS [40] algorithm used in the ARC middleware Peer-to-Peer (P2P) information system. Unlike the basic structured P2P concepts of distributing the keys on an overlay network, and non-structured approaches of replicating the information [28], EMIR slightly modifies the algorithm and replicates the keys among the peer GSR nodes in the network and makes the information eventually consistent [51] after a certain period of time. By replicating the information, all the services can be discovered from any of the available GSRs, which makes the infrastructure resilient to bottlenecks and failures.

The *sparsity*, the number of neighbours each P2P node should replicate to, is another key factor (Fig. 7). Selecting a smaller value would consume less bandwidth at a given time but take longer to reach consistency.



**Fig. 7** EMIR P2P network of registries

#### 4.3.5 Authentication and Authorisation

The DSR and GSR nodes expose a programmatic interface to the service publishers, as well as to the applications, to publish and query service records. In addition, the nodes must connect with the other nodes to form a hierarchy or a P2P network. Publishing a service requires a high LoA credential (X.509 certificate), so attackers can not inject malicious services or modify existing services, so all EMIR nodes, and all entities authorised to publish services, must have X.509 certificates issued by a trusted authority.

#### 4.4 Overall Architecture

Figure 8 is an updated version of Fig. 1, showing the details of the two middle rectangles. To look at this process in more detail, we return to our example from Sect. 4.1.1. Figure 9 depicts a sequence diagram, with the following steps:

1. A CLARIN user requires a EUDAT data sharing service to deposit her data, and therefore send queries for the “data sharing” service types to EMIR.
2. The user sends a request of depositing their research dataset on B2SHARE.
3. As the access token is missing from the user’s request, B2SHARE will redirect (using the HTTP protocol) the user to the B2ACCESS service, the authenticating party, and then further to the user’s organisation IdP.
4. The user authenticates themselves, here with a username and password, to the IdP.
5. We assume for this use case that the user is already registered with the B2ACCESS service, so B2ACCESS will not attempt to register them. Instead, B2ACCESS updates its information about the user, if necessary, based on the user attributes in the SAML assertion which has been received from the IdP.

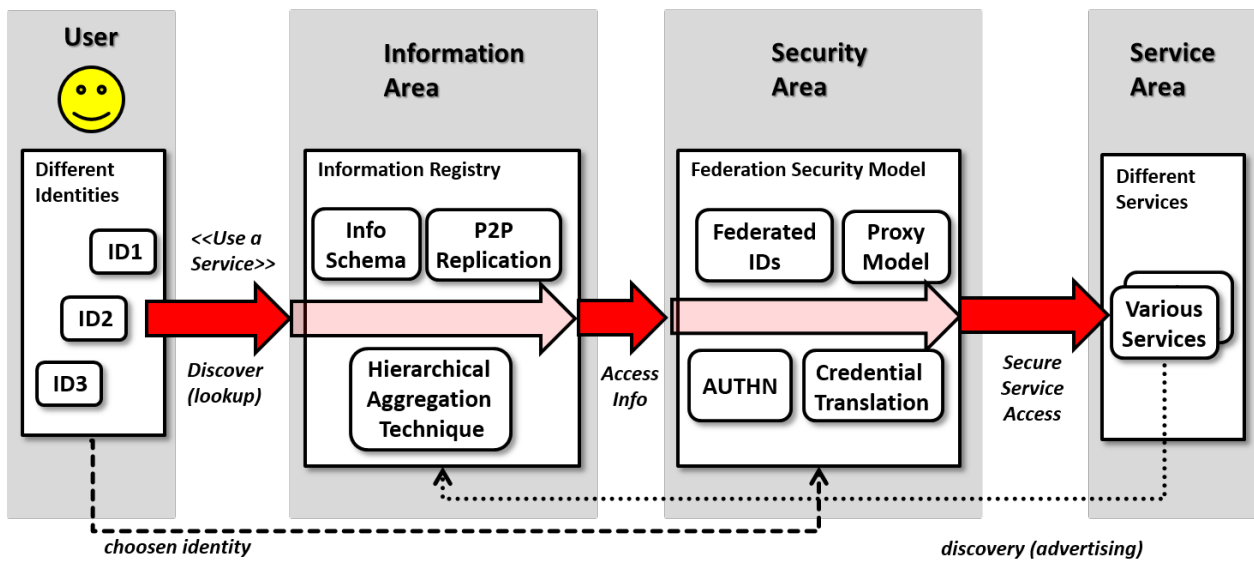


Fig. 8 Integrated federated authentication and service discovery architecture

6. The user (or rather the user's browser) receives and then forwards an authorisation code to the B2SHARE service. On the basis of the code, B2SHARE requests an access token from B2ACCESS.
7. B2SHARE receives an access token.
8. B2SHARE validates the access token and eventually grants the user to deposit/publish/share her data. The data stored on the EUDAT resources (B2SHARE) should now be replicated across multiple storage systems.
9. In order to replicate data with B2SAFE, B2SHARE requires a X.509 credential and sends a request and the access token (from previous flow) to the B2ACCESS Certification Authority (CA) server. The CA server validates the access token and the request.
10. The CA requests a full set of attributes (containing the user's role, group, email, etc.) from the B2ACCESS database.
11. A short-lived X.509 credential is generated, containing the user's attributes in its extensions, and returned to the B2SHARE service.
12. The B2SHARE service can now replicate the data to the relevant B2SAFE nodes.

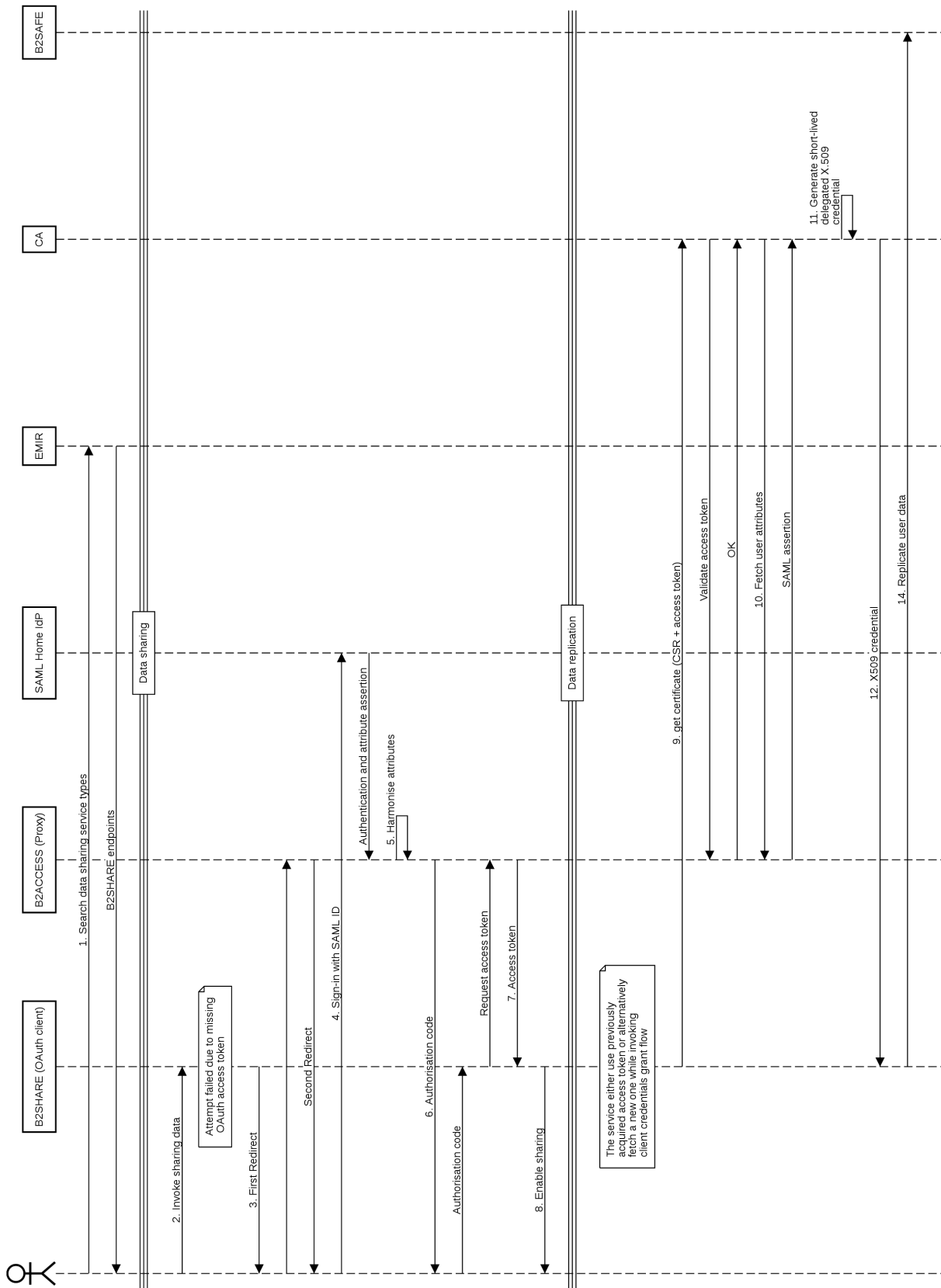
#### 4.4.1 Cross-Infrastructure Federated Service Access

Figure 10 extends the example in Sect. 4.4. We should point out that the scenario is not possible today; a few components are still missing. Nevertheless, it is instructive, as the missing pieces will help us understand the barriers to interoperation.

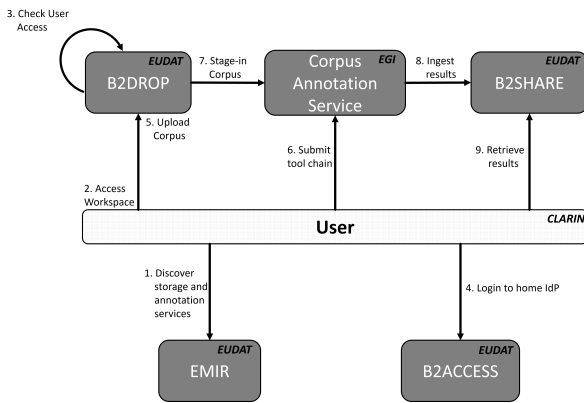
Let us assume that a CLARIN user is in possession of a corpus, and wishes to work on a particular data set from it, consisting of video and image data, and the work will result in annotations.

1. The user looks up data exchange services and corpus annotation services on EMIR, and EMIR returns a list of endpoints on multiple infrastructures.
2. The user selects a B2DROP endpoint (an EUDAT data exchange service [4]) and tries to access its workspace to upload the corpus.
3. B2DROP checks whether the user is authenticated and redirects to B2ACCESS service as before.
4. After successful authentication to B2SAFE, the user uploads the corpus and obtains a unique reference to the corpus. Upload is completed through an OAuth token.
5. The user selects Corpus Annotation Service (CAS) from EGI, and submits a compute request which includes the reference to the corpus, as well as a (bearer) token that authorises the service to access the relevant part of the corpus. The user authenticates to EGI using their EUDAT certificate<sup>1</sup>.
6. CAS retrieves (through the provided reference and token) the corpus and processes the relevant parts of it.
7. B2SHARE receives the processed (annotated) data from the CAS service. Here, B2SHARE does not have any prior authorisation from the user, nor does it have the option to ask for one (as CAS is running without the user's direct intervention). Thus, CAS needs to upload data using the delegated certificate.

<sup>1</sup> If B2DROP had used certificates, the EGI service could have used its delegated certificate to access the data.



**Fig. 9** Sequence diagram showing service discovery, federated authentication, credential translation, and attribute harmonisation in the EUDAT infrastructure



**Fig. 10** CLARIN data staging use case showing cross-infrastructure federated authentication and service discovery

Note that the certificate also contains the e-mail address as metadata, so the service is able to notify the user of the upload, including of course a link to the data.

8. Finally, the user fetches the annotated corpus from the provided link.

As we have mentioned, unlike the scenario in Sect. 4.4, the scenario above is an aim. It is not possible today, but it serves to highlight the current gaps:

- EUDAT and EGI must both accept the same certificates. Unlike IGTF certificates ([www.igtf.net](http://www.igtf.net)), certificates generated for federation-internal use are not trusted across infrastructures, due to the variation in LoA. Work is in progress to harmonise on RCaath [43].
- Likewise, we have, in this scenario, skipped lightly over the authorisation process. In practice, EGI would allocate resources to the community, thus requiring community membership attributes to be communicated *with the credential* because an EGI service would not *a priori* be authorised to query user attributes from B2ACCESS. In fact, these attributes are currently communicated with the credential, but will not be after a migration to RCaath. In other words, cross-infrastructure authorisation needs a lot more thought.
- B2ACCESS provides consistent user mapping across OAuth/OpenID Connect credentials and certificates. In the scenario above, a service would sometimes need to use one, sometimes the other: B2SHARE would need to accept a certificate from CAS, but OpenID Connect from the user’s browser in step 8. In the current infrastructures, services either use one or the other, but not both.
- The current production instance of B2DROP is, as of this writing, not integrated with B2ACCESS.

- An EMIR service is needed which aggregates services across all three infrastructures. Note that there is no access control on querying service information.
- As with resource allocation, accounting also needs to be consistent.

## 5 Discussion

This paper presented a federated AAI and service discovery framework. The B2ACCESS service implementing the AAI presented in this paper fulfils requirements R1, R4, R2, R7, R6, R8, and R9 of CLARIN and EUDAT because it manages authentication, user attributes, and credential translation in one service. In addition to that, EMIR addresses requirements R3, R15, R16, R18, R19, and R20 by offering a robust service discovery for EGI infrastructure (or alike). In particular, it combines a hierarchical model that allows subdomains to manage their resources with a peer-to-peer model across the top-level nodes.

In the context of EMIR, the registry nodes are relatively static in nature, so they can rely on PKI and Access Control Lists (ACLs) for authentication and authorisation, respectively. This requires a communication between the administrators to exchange the nodes’ information.

In terms of B2ACCESS, there are a number of areas (while liaising with EUDAT and AARC) in the future to look into:

- Connecting infrastructures through shared (mutually trusted) authentication. Harmonised communication of the LoA will be useful, which is work in progress through the REFEDS work.
- Supporting multiple LoA and also providing a standard means (e.g. step-up authentication) to augment the assurance levels. Work in progress in AARC should provide guidance on this.
- Integration of a fine-grained and externalised authorisation system based on the XACML standard. However, as we saw in Sect. 4.4.1, much more research is needed in cross-infrastructure authorisation.
- Unsurprisingly, heterogeneous services which need several different “flavours” of credentials (as in EUDAT) make it harder to build cross-infrastructure (or indeed inter-infrastructure) interoperation.

### 5.1 Impact on Infrastructures

B2ACCESS already provides production-ready AAI for EUDAT infrastructure<sup>2</sup>, which implies integration as

<sup>2</sup> <https://b2access.eudat.eu>

well as enabling federated access (using federated identities) to all the B2 services, with dissimilar authentication protocols (SAML, OIDC, PKI). Given the adoption of B2ACCESS in EUDAT, other scientific communities such as EPOS are also considering to deploy B2ACCESS (independently from EUDAT) in their own research infrastructure. B2ACCESS being EUDAT AAI plays an important role within the AARC consortium as one of its objectives is to achieve interoperability of B2ACCESS across e/cyber/research-infrastructures, such as EGI, PRACE and ELIXIR identity and service federations. This is, however, more than an interoperability exercise as (in particular) EGI and EUDAT will have to collaborate by sharing their services within the future EU-funded EOSC-Hub project, the successor of the EUDAT project. Alongside the interoperation, B2ACCESS has fed its experiences into building the AARC Blueprint Architecture [10]. Being an SP/IdP proxy, B2ACCESS has significantly reduced the barrier of trust management between service and identity federations. There are also risks when users' identity is compromised and since EUDAT hosts and manages data from scientific communities, the attacker can delete or rewrite users' datasets with arbitrary data. To cope with such attacks, B2ACCESS adopts the SIRTFI [6] framework to react immediately and mitigate the risks. While the users and services are provisioned into the EUDAT's B2ACCESS service, the registration goes through a formal process for approval by the B2ACCESS administrators, to check whether the identity is compliant with EUDAT policies. As for EMIR, it has been integrated with all the services which are included in the EMI services catalogue, thus it has enabled publishing and querying of the services by the infrastructure operators, monitoring systems and other services (for example workflow). However, EMIR is also being evaluated for service discovery purposes within the EUDAT and EGI infrastructures.

## 5.2 Impact on Users

With B2ACCESS in EUDAT, end users from various scientific communities (CLARIN, ELIXIR, DARIAH, TERENO, EPOS, ENES, etc.) possessing a single identity have federated access to the EUDAT's B2 services. The underlying credentials of the users can be SAML ID, Social ID (from Google, Facebook or ORCID) and X.509 certificates. The EUDAT services do not rely on any single authentication protocol, thus B2ACCESS enabled the authentication by translation of credentials. As far as service discovery is concerned, EMIR has facilitated the users and software clients by querying of

infrastructure services based on the service metadata (service type and capabilities).

## 6 Related work

Federated Identity Management (FIM) or AAI in a broader sense has been a challenge for many years [33]; though the social, commercial and research application providers are recently getting more traction towards external rather than built-in identity management solutions. It is also pertinent for a collaborative infrastructure like EUDAT, providing secure and federated data management services to the research communities [5, 12, 15, 35] in which the earth scientists or linguists would want to collaborate (for example share their data) within or across research communities, given each communities have already their established external or internal identity management system in place, so they bring their own identities.

ELIXIR is one of the largest research infrastructures in Europe, having their own data and identity management infrastructure. The main goals of ELIXIR are to orchestrate the collection, quality control and archiving of large amounts of biological data produced by life science experiments. It also has an aim to improve the long-term sustainability of biological datasets [23]. The ELIXIR AAI [24] provides web identity federation while integrating with Global Authentication Infrastructure (eduGAIN) [20] inter-federation service. In addition to that, the AAI allows users to authenticate with their social identities, which are issued from Google, ORCID and Facebook. It supports associating remote user identities with infrastructure-wide identifiers. Unlike B2ACCESS, ELIXIR AAI lacks support for credentials translation. Similarly, ELIXIR's support for multiple authentication protocols is limited, hence it does not provide end user authentication with end-entity X.509 certificates and LDAP based credentials.

XSEDE [52] is the successor to TeraGrid [36], an NSF funded HPC and grid infrastructure. It consists of a collection of advanced digital resources and services (like supercomputers, visualization and storage systems, collections of data, software, networks, and expert support) that support researchers in various scientific domains. XSEDE relies on Globus Auth [50], a framework for identity and access management. Like B2ACCESS, the Globus Auth framework allows integration with SAML-based identity federations, identity linking, identity brokering (or credential translation) and group management. Furthermore, Globus Auth uses MyProxy-based CILogon [41, 7] to enable federated access to non-browser-based resources, which in particular rely on short-lived X.509 credentials. B2ACCESS

instead uses its own online CA to generate the short-lived credentials. However, integration of B2ACCESS with RCAuth [43] (a modified version of CILogon service for European infrastructures) is being tested and evaluated, but will have consequences, as mentioned above.

## 7 Conclusions

In recent years, large-scale infrastructures have substantially evolved where the federated service discovery and access have become increasingly relevant. Users benefit from having a single credential across the whole infrastructure, and benefit further when it is used across multiple infrastructures. With a unified approach to identity management, authentication, authorisation, and accounting, users are able to run workflows and access and store data from one infrastructure to another, thus further enabling user communities and service providers to build more sophisticated services. As with the registry of services for a country, it should be feasible in the near future to extend these into hierarchies of services, similar to the current global grid infrastructures. However, the details matter, and different technologies, varying levels of assurance, different protocols, schemata, conventions and culture can all provide gaps that prevent users from seamlessly interoperating services across infrastructures. However, as we have seen in the present paper, many of the required building blocks are already present, as is the will to interoperate. Also helpful are the harmonisation activities by REFEDS and AARC, and, if needed, the opportunity for standardisation through standards-defining organisations such as DMTF and OGF.

**Acknowledgements** EUDAT2020 is funded by the EU Framework H2020 – DG CONNECT e-Infrastructures, contract no. 654065 – (Part of) the work reported here was made possible by using the CLARIN infrastructure.

## References

1. Authentication and authorisation research consortium. URL <https://aarc-project.eu>. [accessed 19-November-2016]
2. Alfieri, R., Cecchini, R., Ciaschini, V., dell'Agnello, L., Frohner, A., Gianoli, A., Lörentey, K., Spataro, F.: Voms, an authorization system for virtual organizations. In: F.F. Rivera, M. Bubak, A. Gómez-Tato, R. Doallo (eds.) *Grid Computing, First European Across Grids Conference*, Santiago de Compostela, Spain, February 13-14, 2003, Revised Papers, *Lecture Notes in Computer Science*, vol. 2970, pp. 33-40. Springer (2003). DOI 10.1007/978-3-540-24689-3\_5. URL [http://dx.doi.org/10.1007/978-3-540-24689-3\\_5](http://dx.doi.org/10.1007/978-3-540-24689-3_5)
3. Allcock, W., Bresnahan, J., Kettimuthu, R., Link, M., Dumitrescu, C., Raicu, I., Foster, I.: The Globus striped GridFTP framework and server. In: *Proceedings of the 2005 ACM/IEEE Conference on Supercomputing, SC '05*. IEEE Computer Society, Washington, DC, USA (2005). DOI 10.1109/SC.2005.72
4. B2DROP. URL <https://www.eudat.eu/services/b2drop>. [accessed 5-January-2017]
5. Bailo, D., Jeffery, K.G., Spinuso, A., Fiameni, G.: Interoperability oriented architecture: The approach of epos for solid earth e-infrastructures. In: *2015 IEEE 11th International Conference on e-Science*, pp. 529-534 (2015). DOI 10.1109/eScience.2015.22
6. Barton, T., Basney, J., Groep, D., Harris, N., Johansson, L., Kelsey, D., Koranda, S., Wartel, R., West, A., Short, H.: A security incident response trust framework for federated identity (sirtfi). Recommendation Sirtfi-1.0, REFEDS (2015). URL <https://refeds.org/wp-content/uploads/2016/01/Sirtfi-1.0.pdf>
7. Basney, J., Fleury, T., Gaynor, J.: Cilogon: A federated x.509 certification authority for cyberinfrastructure logon. *Concurrency and Computation: Practice and Experience* **26**(13), 2225-2239 (2014). DOI 10.1002/cpe.3265. URL <http://dx.doi.org/10.1002/cpe.3265>. CPE-13-0334.R1
8. Baur, T., Breu, R., Kálmán, T., Lindinger, T., Milbert, A., Poghosyan, G., Reiser, H., Romberg, M.: An interoperable grid information system for integrated resource monitoring based on virtual organizations. *Journal of Grid Computing* **7**(3), 319-333 (2009). DOI 10.1007/s10723-009-9134-3. URL <https://doi.org/10.1007/s10723-009-9134-3>
9. Grid information system. URL <http://gridinfo.web.cern.ch>. [accessed 5-September-2017]
10. Biancini, A., Florio, L., Haase, M., Hardt, M., Jankowski, M., Jensen, J., Kanellopoulos, C., Liampotis, N., Lichehammer, S., Memon, S., van Dijk, N., Paetow, S., Prochazka, M., Sallé, M., Solagna, P., Stevanovic, U., Vaghetti, D.: AARC: first draft of the blueprint architecture for authentication and authorisation infrastructures. *CoRR* **abs/1611.07832** (2016). URL <http://arxiv.org/abs/1611.07832>
11. Blumtritt, J., Elbers, W., Goosen, T., Hinrichs, M., Qiu, W., Sallé, M., Windhouwer, M.: User delegation in the CLARIN infrastructure. In: *Selected Papers from the CLARIN 2014 Conference*, October 24-25, 2014, Soesterberg, The Netherlands. Linköping University Electronic Press, Linköping, Sweden (2015). URL <http://www.ep.liu.se/ecp/article.asp?issue=116&volume=&article=002>
12. Bogen, H.: Tereno: German network of terrestrial environmental observatories. *Journal of large-scale research facilities* **Vol 2**, A52 (2016). DOI <http://dx.doi.org/10.17815/jlsrf-2-98>. URL <http://jlsrf.org/index.php/lfs/article/view/98>
13. Chadwick, D.W., Siu, K., Lee, C., Fouillat, Y., Geronville, D.: Adding federated identity management to openstack. *Journal of Grid Computing* **12**(1), 3-27 (2014). DOI 10.1007/s10723-013-9283-2. URL <https://doi.org/10.1007/s10723-013-9283-2>
14. Christos, K., Nicolas, L., van Dijk Niels, Peter, S.: Deliverable dJRA1.1: Analysis of user community and service provider requirements. Project Deliverable AARC-DJRA1.1, AARC Project (2015). URL <https://aarc-project.eu/wp-content/uploads/2015/10/AARC-DJRA1.1.pdf>



15. CLARIN. URL <https://www.clarin.eu>. [accessed 13-July-2017]
16. CLARIN services. URL <https://www.clarin.eu/content/services>. [accessed 5-September-2017]
17. Conway, M., Moore, R., Rajasekar, A., Nief, J.Y.: Demonstration of policy-guided data preservation using iRODS. In: 2011 IEEE International Symposium on Policies for Distributed Systems and Networks (POLICY), pp. 173–174 (2011). DOI 10.1109/POLICY.2011.17
18. Cornwall, L.A., Jensen, J., Kelsey, D.P., Frohner, Á., Kouřil, D., Bonnassieux, F., Nicoud, S., Lörentey, K., Hahkala, J., Silander, M., Cecchini, R., Ciaschini, V., dell’Agnello, L., Spataro, F., O’Callaghan, D., Mulmo, O., Volpato, G.L., Groep, D., Steenbakkers, M., McNab, A.: Authentication and authorization mechanisms for multi-domain grid environments. *Journal of Grid Computing* **2**(4), 301–311 (2004). DOI 10.1007/s10723-004-8182-y. URL <https://doi.org/10.1007/s10723-004-8182-y>
19. Drollette, D.: Standards are the glue 2.0. iSGTW (ScienceNode) (2009). URL <https://sciencenode.org/feature/isgtw-feature-standards-are-glue-20.php>
20. eduGAIN. URL <http://www.edugain.org>. [accessed 10-August-2017]
21. EGI. URL <http://www.egi.eu>. [accessed 5-September-2017]
22. Federated cloud information discovery. URL [https://wiki.egi.eu/wiki/Federated\\_Cloud\\_Information\\_Discovery](https://wiki.egi.eu/wiki/Federated_Cloud_Information_Discovery). [accessed 5-September-2017]
23. ELIXIR. URL <https://www.elixir-europe.org>. [accessed 15-September-2017]
24. ELIXIR AAI documentation. URL <https://www.elixir-europe.org/services/compute/aai>. [accessed 13-September-2017]
25. European Middleware Initiative (EMI). URL <http://www.eu-emi.eu>. [accessed 10-June-2016]
26. EUDAT collaborative data infrastructure. URL <http://www.eudat.eu>. [accessed 2-September-2016]
27. Field, L., Memon, A.S., Márton, I., Szigeti, G.: The EMI registry: Discovering services in a federated world. *Journal of Grid Computing* **12**(1), 29–40 (2014). DOI 10.1007/s10723-013-9284-1. URL <http://dx.doi.org/10.1007/s10723-013-9284-1>
28. Forestiero, A., Mastroianni, C., Spezzano, G.: Building a peer-to-peer information system in grids via self-organizing agents. *Journal of Grid Computing* **6**(2), 125–140 (2008). DOI 10.1007/s10723-007-9062-z. URL <https://doi.org/10.1007/s10723-007-9062-z>
29. Foster, I.: Globus toolkit version 4: Software for service-oriented systems. In: Proceedings of the 2005 IFIP International Conference on Network and Parallel Computing, NPC’05, pp. 2–13. Springer-Verlag, Berlin, Heidelberg (2005). DOI 10.1007/11577188.2
30. Foster, I.: Globus online: Accelerating and democratizing science through cloud-based services. *IEEE Internet Computing* **15**(3), 70–73 (2011). DOI 10.1109/MIC.2011.64. URL <http://dx.doi.org/10.1109/MIC.2011.64>
31. Foster, I., Kesselman, C., Tuecke, S.: The anatomy of the grid: Enabling scalable virtual organizations. *The International Journal of High Performance Computing Applications* **15**(3), 200–222 (2001)
32. Hardt, M., (eds), C.K.: Blueprint architecture. Project deliverable, AARC Project (2017). URL <https://aarc-project.eu/documents/deliverables/>. (to appear)
33. Jensen, J.: Federated identity management challenges. In: 2012 Seventh International Conference on Availability, Reliability and Security, pp. 230–235 (2012). DOI 10.1109/ARES.2012.68
34. Jensen, J., Stevanovic, U., Kakavas, I., Liampotis, N., Haase, M., Gietz, P., Jankowski, M., Reale, M., Mantovani, M.L., Florio, L.: Design for deploying solutions for “guest identities”. Project milestone, AARC Project (2016). URL <https://aarc-project.eu/wp-content/uploads/2016/06/MJRA1.2-Design-for-Deploying-Solutions-for-Guest-Identities.pdf>
35. Joussaume, S., Budich, R.: The Infrastructure Project of the European Network for Earth System Modelling: IS-ENES, pp. 5–9. Springer Berlin Heidelberg, Berlin, Heidelberg (2013). DOI 10.1007/978-3-642-36597-3\_2. URL [https://doi.org/10.1007/978-3-642-36597-3\\_2](https://doi.org/10.1007/978-3-642-36597-3_2)
36. Katz, D.S., Callaghan, S., Harkness, R., Pamidighantam, S., Pierce, M., Plale, B., Song, C., Towns, J.: Science on the teragrid **Special Issue** **2010**, 81–97 (2010)
37. Mathieu, G., Richards, D.A., Gordon, D.J., Novales, C.D.C., Colclough, P., Viljoen, M.: Gocdb, a topology repository for a worldwide grid infrastructure. *Journal of Physics: Conference Series* **219**(6), 062,021 (2010). URL <http://stacks.iop.org/1742-6596/219/i=6/a=062021>
38. Memon, A.S., Riedel, M., Field, L., Szigeti, G., Marton, I.: EMIR: An EMI Service Registry for Federated Grid Infrastructures. In: EGI Community Forum 2012 / EMI Second Technical Conference, Munich (Germany), 26 Mar 2012 – 30 Mar 2012, Proceedings of Science. Sissa, Trieste (2012). URL <http://pos.sissa.it/archive/conferences/162/073/EGICF12-EMITC2-073.pdf>
39. MongoDB for GIANT Ideas. URL <https://www.mongodb.com>. [accessed 5-September-2017]
40. NorduGrid: ARC peer-to-peer information system. Documentation and developer’s guide NORDUGRID-TECH-21, NorduGrid (2013). URL [http://www.nordugrid.org/documents/infosys\\_technical.pdf](http://www.nordugrid.org/documents/infosys_technical.pdf)
41. Novotny, J., Tuecke, S., Welch, V.: An online credential repository for the grid: Myproxy. In: Proceedings 10th IEEE International Symposium on High Performance Distributed Computing, pp. 104–111 (2001). DOI 10.1109/HPDC.2001.945181
42. Parducci, B., Lockhart, H., Rissanen, E.: extensible access control markup language (XACML) version 3.0. OASIS Standard xacml-3.0-core-spec-en, OASIS (2013). URL <http://docs.oasis-open.org/xacml/3.0/xacml-3.0-core-spec-en.pdf>
43. Research and Collaboration Authentication Certification Authority Service. URL <https://www.rcauth.eu>. [accessed 16-September-2017]
44. van Rijn, A., Vandenbroucke, R.: Guide to e-infrastructure requirements for european research infrastructures. ISBN 978-90-823661-5-0, E-IRG (2017). URL <http://e-irg.eu/catalogue/eirg-1004>
45. Robertson, L.: Computing Services for LHC: From Clusters to Grids, pp. 69–89. Springer Berlin Heidelberg, Berlin, Heidelberg (2012). DOI 10.1007/978-3-642-23157-5\_3. URL [https://doi.org/10.1007/978-3-642-23157-5\\_3](https://doi.org/10.1007/978-3-642-23157-5_3)
46. Rowstron, A.I.T., Druschel, P.: Pastry: Scalable, decentralized object location, and routing for large-scale peer-to-peer systems. In: Proceedings of the IFIP/ACM International Conference on Distributed Systems Platforms, Middleware ’01, pp. 329–350. Springer-Verlag, London, UK (2001). URL <http://dl.acm.org/citation.cfm?id=646591.697650>

47. Sergio, A., Burke, S., Ehm, F., Field, L., Galang, G., Konya, B., Litmaath, M., Millar, P., Navarro, J.P.: GLUE specification v. 2.0. Recommendation GFD-R-P.147, Open Grid Forum (2009). URL <https://www.ogf.org/documents/GFD.147.pdf>
48. Sergio, A., Burke, S., Field, L., Konya, B., Memon, A.S., Meredith, D., Navarro, J.P., Paganelli, F., Smith, W.: GLUE v, 2.0 – reference realisation to XML schema. Recommendation GFD.209, Open Grid Forum (2013). URL <https://www.ogf.org/documents/GFD.209.pdf>
49. Smith, W., Meredith, D., Memon, A.S., Navarro, J.P.: GLUE v, 2.0 – reference realisation to JSON schema. Recommendation GFD-RP.219, Open Grid Forum (2015). URL <https://www.ogf.org/documents/GFD.219.pdf>
50. Tuecke, S., Ananthakrishnan, R., Chard, K., Lidman, M., McCollam, B., Rosen, S., Foster, I.: Globus auth: A research identity and access management platform. In: 2016 IEEE 12th International Conference on e-Science (e-Science), pp. 203–212 (2016). DOI 10.1109/eScience.2016.7870901
51. Vogels, W.: Eventually consistent. *Commun. ACM* **52**(1), 40–44 (2009). DOI 10.1145/1435417.1435432. URL <http://doi.acm.org/10.1145/1435417.1435432>
52. XSEDE. URL <https://www.xsede.org>. [accessed 13-September-2017]
53. Zinn, C., Hinrichs, M., Dima, E., van Uytvanck, D.: CLARIN switchboard specification. CE-2015-0684, CLARIN (2015). URL [https://office.clarin.eu/v/CE-2015-0684-LR\\_switchboard\\_spec.pdf](https://office.clarin.eu/v/CE-2015-0684-LR_switchboard_spec.pdf)