# VRE4EIC

**A Europe-wide Interoperable Virtual Research Environment
to Empower Multidisciplinary Research Communities
and Accelerate Innovation and Collaboration**

# Deliverable D5.3

# A strategy for the VRE4EIC project to handle security, privacy and trust issues – second version

Document version: 1.0

# VRE4EIC DELIVERABLE

Name, title and organisation of the scientific representative of the project's coordinator:

Mr Philippe Rohou      t: +33 4 97 15 53 06       f: +33 4 92 38 78 22       e: philippe.rohou@ercim.eu

GEIE ERCIM, 2004, route des Lucioles, Sophia Antipolis, F-06410 Biot, France

Project website address: http://www.vre4eic.eu/

| Project | |
|---|---|
| Grant Agreement number | 676247 |
| Project acronym: | VRE4EIC |
| Project title: | A Europe-wide Interoperable Virtual Research Environment to Empower Multidisciplinary Research Communities and Accelerate Innovation and Collaboration |
| Funding Scheme: | Research & Innovation Action (RIA) |
| Date of latest version of DoW against which the assessment will be made: | 31 May 2017 Amended Grant Agreement through amendment n°AMD-676247-8 |
| **Document** | |
| Period covered: | M1-M28 |
| Deliverable number: | D5.3 |
| Deliverable title | Strategy for the VRE4EIC project to handle security, privacy and trust issues - second version |
| Contractual Date of Delivery: | 31/01/2018 |
| Actual Date of Delivery: | 30/01/2018 |
| Editor (s): | Laura Hollink |
| Author (s): | Daniele Bailo, Valerie Brasse, Tessel Bogaard, Cesare Concordia, Laura Hollink, Paul Martin, Jacco van Ossenbruggen, Jan Wielemaker, Yi Yin |
| Reviewer (s): | Maria Theodoridou, Carlo Meghini |
| Participant(s): | All |
| Work package no.: | 5 |
| Work package title: | Information management policy, security, privacy and VRE trustability |
| Work package leader: | Laura Hollink |
| Distribution: | PU |
| Version/Revision: | 1.0 |
| Draft/Final: | Final |
| Total number of pages (including cover): | 45 |

# What is VRE4EIC?

VRE4EIC develops a reference architecture and software components for VREs (Virtual Research Environments). This e-VRE bridges across existing e-RIs (e-Research Infrastructures) such as EPOS and ENVRI+, both represented in the project, themselves supported by e-Is (e-Infrastructures) such as GEANT, EUDAT, PRACE, EGI, OpenAIRE.  The e-VRE provides a comfortable homogeneous interface for users by virtualising access to the heterogeneous datasets, software services, resources of the e-RIs and also provides collaboration/communication facilities for users to improve research communication.  Finally it provides access to research management /administrative facilities so that the end-user has a complete research environment.

# Disclaimer

This document contains description of the VRE4EIC project work and findings.

The authors of this document have taken any available measure in order for its content to be accurate, consistent and lawful. However, neither the project consortium as a whole nor the individual partners that implicitly or explicitly participated in the creation and publication of this document hold any responsibility for actions that might occur as a result of using its content.

This publication has been produced with the assistance of the European Union. The content of this publication is the sole responsibility of the VRE4EIC consortium and can in no way be taken to reflect the views of the European Union.

The European Union is established in accordance with the Treaty on European Union (Maastricht). There are currently 28 Member States of the Union. It is based on the European Communities and the Member States cooperation in the fields of Common Foreign and Security Policy and Justice and Home Affairs. The five main institutions of the European Union are the European Parliament, the Council of Ministers, the European Commission, the Court of Justice and the Court of Auditors (http://europa.eu/).

VRE4EIC has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 676247.

# Table of Contents

# 1 Introduction

An e-VRE provides researchers across research domains the means to share and reuse data and computations. This goal means an e-VRE has to be highly ambitious regarding the level of security to protect both the shared research data and the personal data of the e-VRE users. Moreover, it requires that highly accurate metadata is available for users to reliably assess how much they can trust the resources (data, software, or other resources) that are made available to them. In this deliverable, we present a strategy for the VRE4EIC architecture and software components regarding trust, privacy and security issues for the perspective of the VRE4EIC project. **This is an updated version, adding new insights and improvements to the first version of the strategies presented in deliverable D5.1.** The same issues from the perspective of the end-user are the topic of deliverables D5.2 and D5.4. Upon acceptance of this strategy document, it will be made publicly available and especially distributed to EPOS and ENVRIplus and other use case partners.

This document covers questions related to the data or other resources that are shared via the e-VRE, as well as questions about how to provide access to those resources. The latter is frequently referred to as AAAI (Authentication, Authorization, Accounting Infrastructure).

## 1.1 Updates with respect to version 1 of this document

Firstly, we have sketched the dynamic, international context in which AAAI solutions operate. Since the delivery of the first version of this document, AAAI solutions have continued to develop and plans for future solutions have crystallized. Many solutions are still unfinished, notably attribute management and the organizational aspects of federated login. The monitoring of these ongoing developments has informed both our long-term recommendations as well as our proposals for practical solutions. Our reflection on and conclusions about these issues can be found in section 3.

Secondly, the Trust strategy has been extended with methods to ensure the integrity of (meta)data (section 4.2).

Thirdly, we have updated the strategy regarding Security with detailed descriptions of the options that exist when implementing two-factor authorization (section 7.4)

Finally, we have extended the Metadata strategy with recommendations for a crucial issue for secure Role-Based Access Control (RBAC) in the e-VRE: who assesses, and guarantees, the quality of the role assignments that are used for access control? (section 8.4)

These updates have lead to four new recommendations including pointers to affected parts of the architecture (namely Trust recommendation 4 and Security recommendations 8, 9 and 10 in section 10).

## 1.2 Security, privacy and trust of datasets in an e-VRE

Secure storage, transport and backup of data is a core service of an e-VRE. This is especially important for privacy-sensitive data, such as medical records or usage logs. The privacy-sensitive nature of data - even of public data or anonymised data - cannot be taken lightly, as two well-known cases have shown. In 2006, AOL released search logs in which individual users could be identified despite

anonymisation efforts[1], leading to a class action lawsuit. In 2014, Danish researchers led by Emil O. W. Kirkegaard published a dataset[2] containing usernames, age, gender, etc. from the dating site OKCupid. The researchers did not anonymise the data stating that the "Data is already public", spurring a huge ethical debate.

Privacy of data is an issue at both the e-RI and the e-VRE level. However, the fact that the e-VRE bridges across several e-RIs poses additional challenges with regard to privacy. Both national and EU legislation are addressing these issues, and different countries might have different, even conflicting, laws. Moreover, combining datasets can in certain situations cause a privacy breach; a combination of a closed dataset with an open dataset where the open dataset becomes identifiable by joining it to the closed dataset (such as combining Accident Data, available both in the UK and the Netherlands, to the customer data of an insurance company); or the publication of the NetFlix dataset in the now infamous NetFlix Prize, where the applied anonymisation did not suffice; or research showing how little mobility information is needed to identify 95% of all people [De Montjoye 2013]. The privacy-policy of an e-VRE will therefore often be stricter than that of an e-RI.

Trust is necessary for a researcher that uses the e-VRE to assess whether the quality of the resources is sufficient for her needs at that time. We distinguish trust in people (e.g. data owners), trust in software, tools, or algorithms (e.g. those that were used to produce the data) and finally trust in the datasets themselves. These three are interconnected, as will be discussed in Section 3. Most examples in this deliverable focus on trust in data, but that cannot be seen as separate from trust in the involved people and tools. The fact that an e-VRE bridges across research communities has implications for the level of trust that a researcher places in a resource that is accessed through the e-VRE. Firstly, users will typically access resources from other communities, where they are not familiar with the standard data collection processes and where they do not know the other actors. Secondly, the e-VRE will typically be used to combine resources from different places. Uncertainty about the quality of the individual datasets may lead to a higher level of uncertainty about the quality of the combined dataset. To increase transparency and reproducibility, a trust-policy of an e-VRE should be focussed on providing the necessary information for a user to determine whether she can trust a particular dataset. For that, an e-VRE relies on the availability and interoperability of metadata from the e-RIs. Rather than making the e-VRE work exclusively with e-RIs that meet high metadata standards, a preferred strategy is to *incentivise* e-RIs to meet these standards, e.g. by offering better visibility and usability of their data and services. Trust is also tied to the level of security and privacy that a research infrastructure offers. These will be discussed as separate issues in this deliverable.

# 1.3 AAAI

Authentication can be seen from a service perspective (is this really the user s/he claims to be?) and a user perspective (is this service really the service it claims to be?). Authentication requirements vary widely and so does the used technology. We find IP-address based access, simple password restricted access, federated identity services based on e.g., SAML[3] (Security Assertion Markup Language), OAuth2[4], OpenID[5] and/or certificate controlled access based on X509[6]. A successful e-VRE is compatible with a wide variety of identity providers in order to suit the needs of associated e-RIs across multiple countries. In this deliverable, we ignore the many eID schemes developed by EU member

---

[1] https://en.wikipedia.org/wiki/AOL_search_data_leak

[2] http://openpsych.net/forum/showthread.php?tid=279

[3] http://docs.oasis-open.org/security/saml/v2.0/saml-conformance-2.0-os.pdf

[4] https://tools.ietf.org/html/rfc6749

[5] http://openid.net/specs/openid-connect-core-1_0-final.html

[6] https://tools.ietf.org/html/rfc5280

countries for e-government services, since these typically only work on a national scale and cannot yet be used to access research infrastructure on a European scale.

For privacy reasons, federated identity providers are typically reluctant in providing attributes about the confirmed identity (name, email and CERIF data about organisations and projects the user is involved in), resulting in additional challenges for the e-VRE when it comes to granting access to e-RI data. In this policy document, we describe these challenges and provide options for how to deal with them. Just like the identity providers, the e-VRE itself should have a policy for how to guarantee the privacy of its users when storing authentication and access logs.

## 1.4 Structure of this document

We start with an inventory of the relevant findings from the requirements analysis and characterization of existing e-RIs in work package 2, leading to the identification of gaps in how the user needs are being met at both the level of the e-RIs and the level of the e-VRE (Section 2). Section 3 discusses the broader context of (inter)national e-VRE AAAI solutions that are currently being used or under development. Sections 4-8 present strategies for how an e-VRE should deal with issues regarding trust, accounting, privacy, security, and finally the role of (CERIF) metadata. In Section 9, these strategies are translated into specific requirements for the overall e-VRE architecture and software components. Finally, Section 10 contains a list of recommendations for the e-VRE architecture components and policy documents.

# 2 Requirements, existing solutions and gaps

In work package 2 (WP2 from now on for brevity), general requirements for an e-VRE were identified by means of a literature study. In addition, five existing e-RIs were characterized using questionnaires to systematically describe the current state of these systems. The results of these studies have in part been reported in D2.1. This is considered a 'living document'; it will be updated as new insights emerge, for example from additional characterizations of e-RIs. Here, we highlight those findings that have direct implications for security, privacy and trust issues. The two studies are complementary: the requirements analysis gives an overview of the functionality that is deemed necessary for researchers collaborating through an e-VRE; the e-RI characterization provides details of the actual implementation of the e-RIs. In the current section, we **list the main findings of the WP2 studies and identify gaps in how the user needs are being met at both the level of the e-RIs and the level of the e-VRE**.

The five currently characterized e-RIs:

1) EURO-ARGO[7] is a distributed RI for environmental science. It works towards providing research data and related services on climate and oceanography.

2) ELIXIR[8] is a distributed RI for life science. It aims at integrating and sustaining bioinformatics resources generated by publicly funded research and sharing access to the data for Europe's life-science research organisations.

3) The Integrated Carbon Observation System (ICOS)[9] is an RI on environmental science providing long-term observations required to understand the present state and predict future behaviour of the global carbon cycle and greenhouse gas emissions and concentrations.

4) EPOS[10] plans to integrate the Research Infrastructures for Solid Earth Science in Europe. Its goal is to establish a comprehensive multidisciplinary research platform for the Earth sciences.

5) LifewatchGreece [11] Research Infrastructure (LWG RI), funded by the GSRT (structural funds), is the national effort to research on biodiversity data and data observatories in Greece.

For clarity, we follow the structure of information in WP2 documents as much as possible: we indicate the original numbering of each requirement as it was used in D2.1, and we maintain the structure of the interview protocol where data was grouped into topics (relevant topics to security, privacy and trust issues are: data access, licensing for data use, data ownership, liability of data use and data disposal.)

One thing that was noted during the requirements elicitation process is that e-RIs should have a research data management policy and that the users of the e-RI should have an agreement on how to use the data (requirement number SRQ8). For this reason, it is important that the outcome of this deliverable is not only seen as a project deliverable and input for the e-VRE architecture and metadata model, but also as a starting point for a research data management policy that users of the e-VRE are aware of, have access to and agree on.

---

[7] http://www.euro-argo.eu/

[8] https://www.elixir-europe.org/

[9] https://www.icos-cp.eu/

[10] https://www.epos-ip.org/

[11] https://www.lifewatchgreece.eu/

## 2.1 Security

The requirements elicitation process has resulted in a number of concrete functionalities that improve the security of the e-VRE (Table 1). There is a clear need for access control (CRQ6), in particular the ability to provide or deny access to resources to specific users or groups of users. In addition, we identified a need for physical access control (SRQ15) in addition to digital access control. Data should be securely stored (SRQ12) and transmitted (CLRQ15). In order to prevent data loss, the e-VRE should facilitate backup and/or curation of datasets (PRQ35) Finally, there is a need for accounting services (PRQ31) , which includes the logging of data access and use (SRQ6).

| No. | Requirement | Description |
|---|---|---|
| CRQ6 | Data Storage & Preservation | Ability to deposit (over long-term) the data and metadata or other supplementary data and methods according to specified policies, and make them accessible on request. |
| SRQ15 | Physical access control | Identity control of the access to the physical infrastructure |
| SRQ12 | Secure storage | Secure storage of data, especially sensitive data |
| CLRQ15 | Data Transmission | Ability to transfer data over communication channel using specified network protocols. |
| PRQ35 | Data backup | Ability to backup datasets according to specified policies |
| SRQ6 | Use log | Logs of the system usage for auditing and legal compliance |
| PRQ31 | Accounting | Accounting services for data and services provider |

Table 1: Identified requirements related to security.

The e-RI characterizations show that a variety of AAAI solutions are currently in place or under development at the e-RIs. Not all data in the e-RIs requires restricted access; some data is open and freely available (Table 2).

| e-RI | Data access |
|---|---|
| **EURO-ARGO** | All ARGO data are publicly accessible. There is no restriction on the use of published data. No login is required. |

| ELIXIR | Open access is provided to all publicly available data; and secure controlled access is provided to sensitive personal data. No login is required to access the publicly available data. |
|--------|-----------------------------------------------------------------------------------|
| ICOS | The basic rule is that the data does require access restrictions. Some premium services however might be only available to people who have signed up for it. Also if a user wants to save previous searches, this is only possible when he/she has a profile. All data products are free. Final data products are available via the ICOS Carbon Portal. Other types/levels of data can be obtained via the Ecosystem Thematic Centres or from the PI of the observation stations. The Carbon Portal provides a single sign-on functionality. |
| EPOS | Login and password access with credentials from various providers. EPOS provides open access to  85 % of its data. Only a small amount of data is not open, either subject to an embargo period (6 months) or paid data.<br><br>An AAAI mechanism (UNITY[12]) is planned to be used. |
| LifewatchGreece | Access control is used in the RI |

**Table 2: e-RI characterization result regarding data access**

At the e-VRE level, this means that the e-VRE should (1) be compatible with several external access mechanisms, (2) be able to include new ones when new e-RIs connect to the e-VRE and (3) allow unrestricted access to open data. In the second case, the e-RIs should be warned for potential additional privacy risks when their data is combined with other datasets (differential privacy).

Physical access control (in addition to the standard digital access control) is not used at the five currently characterized e-RIs, while it was identified as a requirement. Individual e-RIs should determine how much priority this requirement has for their user community.

The logging of user actions and accounting may be implemented at both the e-RI level and the e-VRE level. The e-VRE logs allow for a complete picture of user actions across the various e-RIs. Note that while accounting relies on an identification of users, logging of actions of non-registered users is useful as well to provide overall usage statistics.

Secure data-storage, backup and secure transmission of data are handled at the e-RI level. Here, the task of the e-VRE is to provide (CERIF) metadata about the provided level of security, e.g. whether encryption is used.

---

[12] http://www.unity-idm.eu/

## 2.2 Privacy

The e-VRE should guarantee the privacy of both users of the e-VRE and of sensitive research data that is stored through the e-VRE. Access Control (CRQ6), secure storage (PSRQ2) and transmission (DRQ14) of research data were already mentioned as security-related requirements. We mention them under privacy again since they are fundamental to protecting privacy-sensitive research data. In addition, the identities, access credentials as well as transaction logs of users of the e-VRE should be stored securely (PSRQ3). This includes the metadata stored in CERIF.

**Differential privacy:** The fact that the e-VRE bridges across several e-RIs poses additional challenges with regard to privacy. In D2.1 on requirements elicitation, it was noted that *"Datasets often require removing privacy sensitive variables [...] before publication. [...] Moreover, the combination of data with other sources might still make it possible to track the identity of an individual person, especially when open data are combined with social media data. "* This means that the privacy levels of data in an e-RI are not always strict enough for an e-VRE. This results in additional requirements related to resetting access control settings (e.g. to disallow combination of data when an e-RI becomes part of the e-VRE), creating awareness with data providers (that their previous privacy policy might no longer be enough).

| No. | Requirement | Description |
|---|---|---|
| CRQ6 | Data Storage & Preservation | Ability to deposit (over long-term) the data and metadata or other supplementary data and methods according to specified policies, and make them accessible on request. |
| SRQ12 | Secure storage | Secure storage of data, especially sensitive data |
| CLRQ15 | Data Transmission | Ability to transfer data over communication channel using specified network protocols. |
| SRQ13 | Credentials protection | Ability to  protect the user's' digital identities and credentials |

Table 3: Identified requirements related to privacy.

The e-RI characterization process revealed one additional issue that is strongly related to privacy: the ability to remove data from the e-RI. Three out of the five e-RIs are known to have functionality and/or policies in place to make this possible, for varying reasons (Table 4). ICOS is the only e-RI where the possibility to remove personal data is made explicit, even though 'the right to be forgotten' is an EU-wide law.

| e-RI | Data disposal |
|---|---|
| **EURO-ARGO** | When a float (sensor object) enters a certain country's territorial waters, the data transmission can be stopped upon the request from the related country. |

| ELIXIR | Retraction of data is possible for reasons of copyright infringement or personal security (ethical issues). |
|---|---|
| ICOS | User profiles and associated user information can be deleted upon request. |
| EPOS | Unknown |
| LifewatchGreece | Data cannot be deleted upon the request of users. |

**Table 4: e-RI characterization result regarding data disposal.**

None of the characterized e-RIs have mechanisms in place to deal with the specific privacy problems that arise when combining datasets.

## 2.3 Trust

The elicited requirements show a clear need of users for methods to cite data (IRQ4): they need to be able to uniquely identify datasets (IRQ1), including parts of datasets (IRQ1) or specific versions of datasets (CRQ4); they need a guarantee that identified data will not change and remain accessible (CRQ6). This enhances the reproducibility of studies done on the basis of these data. In addition, these identification mechanisms provide a means to keep track of changes made to datasets, in other words, to record the provenance. Finally, the opportunity to verify the quality of the data (CRQ6 and CRQ3), improves the transparency of the research process.

We observe that in some cases there may be a tension between the need to record provenance of datasets, including information on who did what, and the need to protect the privacy of users, including their identities and access logs (SRQ6 in Section 2.1 above). An e-VRE needs to have a clear policy regarding this issue.

| No. | Requirement | Description |
|---|---|---|
| IRQ1 | Data Identification | Ability to assign (global) unique identifiers to data contents. |
| CRQ4 | Data Versioning | Ability to assign a new version to each state change of data, allow to add and update some metadata descriptions for each version, and allow to select, access or delete a version of data. |

| CRQ6 | Data Storage & Preservation | Ability to deposit (over long-term) the data and metadata or other supplementary data and methods according to specified policies, and make them accessible on request. |
|---|---|---|
| CRQ6 | Data Quality Checking | Ability to detect and correct (or remove) corrupt, inconsistent or inaccurate records from data sets. |
| CRQ3 | Data Quality Verification | Ability to support manual quality checking. |
| CRQ7 | Data Replication | Ability to create, delete and maintain the consistency of copies of a data set on multiple storage devices. |
| CLRQ18 | Data Publication | Ability to provide clean, well-annotated, anonymity-preserving datasets in a suitable format, and by following specified data-publication and sharing policies to make the datasets publicly accessible or to those who agree to certain conditions of use, and to individuals who meet certain professional criteria. |
| IRQ4 | Data Citation | Ability to assign an accurate, consistent and standardised reference to a data object, which can be cited in scientific publications. |

**Table 5: Identified requirements related to trust.**

From the characterizations it became clear that the e-RIs spend considerable effort on making explicit who owns a dataset (Table 7) and what one is allowed to do with the dataset (licensing and liability in Table 6 and Table 8 resp.) This information is provided e-RI-wide or as metadata with the individual datasets.

| e-RI | Licensing for data use |
|---|---|
| **EURO-ARGO** | CC BY 4.0 |
| **ELIXIR** | Licence agreements are currently under expansion, revision and discussion in ELIXIR |
| **ICOS** | CC BY 4.0 |
| **EPOS** | CC BY 4.0 |

| LifewatchGreece | Data providers choose a licence and embargo period during the upload of the data |

Table 6: e-RI characterization result regarding licensing.

| e-RI | Data ownership |
| --- | --- |
| EURO-ARGO | The data belongs to the owner of the sensor object that collects the data (the 'float'). Upon joining the EURO-ARGO initiative, the one who purchased the float needs to agree that all data will be made freely accessible. |
| ELIXIR | There are policies regarding the data ownership, and metadata about data ownership is provided. |
| ICOS | There is a policy regarding data ownership and there will be metadata about data ownership, including developer, contributor, institution, funding, contact and citation. |
| EPOS | There is metadata about data ownership. |
| LifewatchGreece | There is metadata about data ownership |

Table 7: e-RI characterization result regarding data ownership.

| e-RI | Liability of data use |
| --- | --- |
| EURO-ARGO | There are disclaimer terms in the user agreement. |
| ELIXIR | There are policies regarding the liability of data use. |
| ICOS | A data policy is under development. |
| EPOS | There is data management policy about this. |
| LifewatchGreece | No information regarding this topic. |

Table 8: e-RI characterization result regarding liability.

At the e-VRE level, the main requirement is to correctly convey the information that is already present at the e-RI level (incl. data ownership, licensing and liability) of each dataset as metadata, preferably in CERIF.

# 3 Organizational and technical developments

Our strategy should be aligned with international developments regarding AAAI.  This section contains insights from literature study and discussions with representatives of the Dutch SURFsara and CLARIAH e-VRE project about current and future IAAA solutions. The situation may be slightly different in other countries, but the overall direction will be the same, as it is dictated by standards and new software components that are developed in international contexts.

An e-VRE should enable single sign-on based on the identity provided by the involved research institutes.  This technique is known as federated identity.  This field was pioneered by OpenID and is now common practice on the web. A typical example is a site that allows users to login with their Google or Facebook account. The basic technology is provided by the SAML and OAuth2 protocols discussed above. These protocols allow a relying party (a site to which the user wants access) to delegate the login to the identity provider (e.g., Google).  The identity provider performs the login and confirms to the relying party that this user is known. The identity provider provides the relying party with an access token, identity (a token) and optionally attributes such as the user's email and name.

This technology is at the core of the eduGAIN initiative to provide access to research infrastructures similarly as researchers now have access to the wifi network of other research institutes using eduroam.  The initiative faces a number of both technical and organizations challenges.  Below we describe these with examples from the Netherlands.

- The SAML and oauth2 protocols are designed for logging in to web sites. In this scenario the identity of the user is established and a *session cookie* is used to maintain the authenticated state of the user. Research environments, however, often use different authentication protocols such as X.509 certificate based protocols (e.g., ssh) and proprietary protocols used by databases. Initiatives such as COmanage[13] provide an infrastructure to associate an X.509 *public key* with a federated identity.   This, however, is not yet established technology.
- The SAML and oauth2 protocols require the *relying party* (see above) to be registered with the identity provider.  This is needed to guarantee the security of the protocol.  At the same time this can be used as a statement of trust that the identity provider considers the relying party trustworthy. Note that this is not required.  For example, anyone can create a Google account and use this to register any website for 'login with Google'.  Google does not verify that the website is trustworthy.  On the other hand, the Dutch DigiD[14] that can be used to login with government agencies and some private organizations such as insurance companies has a very strict procedure for accepting new relying parties. This may be compared to TLS certificates that may be obtained from Let's Encrypt[15] that only verifies the DNS really belongs to the site and the requester of the certificate can indeed control the contents of the site or one of the traditional certificate providers that also validate that the requesting organization exists and is trustworthy.
- The protocols allow providing *attributes* of the user known to the identity provider to the relying party.   Normally the relying party requests a *profile* about the user.  On the first login the identity provider confirms with the user that, for example, his or her email address is requested and whether or not this is acceptable.   While attributes such as name, email,

---

[13] https://spaces.internet2.edu/display/COmanage/Home

[14] https://www.digid.nl/

[15] https://letsencrypt.org/

scientific role and detailed affiliation are in general required for an RI to grant access, such attributes are rarely available from the identity provider for privacy reasons.

In practice, academic and research institute's identity services have strict rules to accept relying parties for their services and are much more reluctant than e.g., social login services such as Google and Facebook in providing attributes that reveal the true identity of the identified person. For example, academic employees and students in the Netherlands can buy software at reduced prices from SURFspot[16]. Users can prove they are entitled to do so using a federated login using their institutes identity provider. The only attribute provided to SURFspot, however, is the fact that the user is entitled to buy software at their site.

The practice described above clearly does not support the requirements for research environments: it is practically impossible to register with all institutions from which the RI may wish to grant users access and there is insufficient data about new users to decide whether or not to grant access to the RI. Even for e-VREs such as the Dutch Clariah[17] project providing a research infrastructure to humanity scholars, this proved impossible from an organizational perspective.

A proper federated login infrastructure for research environments requires a single broker that acts as an intermediate between research institutes and all research environments. Such initiatives are starting at the national level, for example Science Collaboration Zone (SCZ)[18] is an attempt to create such an entity for the Netherlands. The technical infrastructure acts as a SAML/oauth2 *proxy*, relying on the research institutes and providing the required attribute management and support for e.g., X.509 certificates. The scenario is as follows:

- A research institute wishing to use this identity service creates a profile that includes the required attributes about users and the SAML/oauth addresses.
- A person willing to use this service is redirected to SCZ, which on first access provides a list of cooperating institutions. The user selects his or her institute and identifies his or herself. The SCZ asks the user to fill out all attributes that are not already known and required by the RI.
- On completion, a request for access is sent to the RI administrative staff. The staff may accept or reject the request, possibly after communication with the applicant institute.
- After acceptance, the user can login using the SCZ federated identity service and may register X.509 certificates when required by the RI.

From a practical viewpoint we reach the following conclusions:

- Federated login using research institutes identity services needs to be agreed upon at another organizational level than e-VREs. It currently takes place at the national level and eventually shall take place at European or world level.
- Until such services are available we should adhere to a local replacement that uses the same principles such that RIs can easily switch to the new services when they become available. Our current solution is described in section 7.
- Identification will be based on HTTP redirect and SAML/oauth2 or a successor.

---

[16] https://www.surfspot.nl/

[17] https://www.clariah.nl/

[18] https://wiki.surfnet.nl/display/SCZ

● Such services will deal with attribute management.  There is no established standard on how these attributes are defined, entered and verified. In section 8, we discuss how this would work in combination with CERIF. Other systems represent them as a flat set of name/value pairs.

# 4 Trust strategy

## 4.1 Trust in the cross-RI environment of the e-VRE

Various definitions of trust have been proposed [Artz 2007, Golbeck 2006, Ceolin 2014]. Castelfranchi and Falcone [1998] (as described by Sabatter and Sierra [2005]) formulate it as follows: "the decision that an agent X (trustor) takes to delegate a task to agent Y (trustee) is based on a specific set of beliefs and goals, and this mental state is what we call trust." In the case of the e-VRE, the trustor is a user that accesses the e-VRE to fulfill a goal; the trustee is either a user that provides datasets through the e-VRE or a system/algorithm that was used to produce or process these datasets. The trust that a user places in a dataset is directly related to the trust she places in the provider of that data or the systems used to produce and process that data. The same holds for other types of resources shared through the e-VRE, such as tools, code or articles.

We expect that the a-priori level of trust is lower when working through an e-VRE than an e-RI. Firstly, in an e-VRE, users will typically access resources from other communities, where they are not familiar with the standard data collection processes and where they don't know the other actors. Secondly, the e-VRE will typically be used to combine resources from different places. Uncertainty about the quality of the individual datasets may lead to a higher level of uncertainty about the quality of the combined dataset. In general, the VRE4EIC strategy for trust issues is based on the concepts of transparency [Stevens 2000], reproducibility [Bechhofer 2013] and empowerment of the e-VRE user. The e-VRE has been designed such that the user is provided with the necessary information (from the CERIF metadata model) so that she can:

A. Assess whether the quality of the resources is sufficient for her needs at that time. This may require users to be trained in order to be able to explicitly formulate the required quality for their research goals.
B. Cite the resources that she used so that her work becomes transparent (e.g. discovered errors may be traced back to errors in the underlying datasets) and reproducible.

The e-RIs may provide several types of information that contribute to an assessment of how much a dataset or a system/algorithm can be trusted:

1. **Provenance information.** The trustworthiness of a resource can only be determined if information about its creation is available. This typically includes information about the actors, tools and processes involved in producing it or modifying it, including the chronology of the events. Standardization is especially important in an e-VRE context where provenance information needs to be shared across platforms. Several standard formats for provenance information exist, the most notable being PROV-O[19], a W3C recommendation. Provenance information can also be encoded in CERIF, and mappings between CERIF and PROV-O exist [Compton 2014]
2. **Resource ownership information.** Trust is increased by the availability of a name and contact details of the person/organisation that can be held responsible for the quality of the resource, and that is available for questions about the resource. This information may be captured as part of provenance information.
3. **Versioning information.** At least, this includes the relations between different versions of resources, i.e. the fact that one dataset is a newer version of another dataset. In that case, the versioning information can be represented in a provenance metadata format. A more elaborate versioning system might also include information on how and which part of a

---

resource has changed in a newer version. For example, which items in a dataset were merged or deleted in a newer version (by what software/person/organization and for what purpose).

4. **Connections between (raw) data and data products.** Data products such as samples and processed versions should be traced back to the original dataset they were derived from.

5. **Explicit quality information.** Two out of five e-RIs that have been characterized offer the means to (manually) check in advance the quality of data, and to store the outcomes of these checks in metadata fields. For example, one can check that data values fall into a realistic range to assess the proper operation of a remote sensor.

6. **Certificate information.** Whether datasets or processing procedures have been verified by a specific industrial standard can influence the trust. Certification can be 1) basic certification - Data Seal of Approval (DSA), 2) extensive certification - DIN 31644(7) standard: 'Information and documentation - Criteria for trustworthy digital archives', and 3) formal certification - ISO 16363(8).

7. **Open source code.** This gives the researchers a way to verify technical solutions.

**The strategy of the e-VRE will be to provide functionalities for users to access the trust-related information of the e-RIs** when working with a resource via the e-VRE. This strategy is preferred over the alternative - to include functionality to record and store trust-related information within the e-VRE itself. It avoids situations where provenance/versioning/quality information for a resource exists in different, possibly disconnected places, leading to potentially inconsistent and incomplete information seen by a user.

This strategy requires **interoperability between the metadata formats in use at the underlying e-RIs**. This is typically realized through a shared metadata model, a role that is taken by CERIF (WP4). The more interoperable the metadata formats are, the better the e-VRE can provide a user with consistent trust-related information, allowing, for example, a comparison of datasets.

A requirement at the side of the e-RIs is to provide unique identifiers of datasets, parts of datasets, and versions of datasets, and to guarantee the permanence of these resources. If they do not meet these requirements, the trust that a user can have in a dataset (or in the process that includes the dataset) is going to be limited. Again, **the strategy of the e-VRE will be to present a user with clear information about the permanence of a resource**. For example, when a user designs a workflow that includes a dataset for which the e-RI cannot guarantee its stability, this should be flagged as such. Similarly, the search and selection functionality of the e-VRE (e.g. used to list resources in all e-RIs that meet certain criteria) should allow users to limit search results to properly identified, permanent resources.

The above are examples of our **general strategy to incentivise e-RIs to implement advanced trust functionalities and to publish interoperable metadata** about their resources by offering better visibility and usability of their resources. There is a tradeoff between on the one hand the need to record who did what with which resource and on the other hand the privacy of e-VRE users. In this case, we use the strategy of incentives described in Section 1, operationalized in two guidelines: (1) a user's identity is never made public unless she voluntarily agrees to this; (2) we use incentives to convince people to assert ownership of their resources so they can be acknowledged appropriately. The incentives include more visibility of their work and better traceability of their resources, leading to increased trust of other users in their resources, and therefore to more reuse and ultimately a higher citation count. To strengthen this incentive, the e-VRE aims to make it easy for people to cite datasets and software/algorithms (in cases where these resources have a permanent identifier). The implementation details of this aim are to be decided at a later stage. One solution would be to include an 'export citation' option, as is currently common practice in online publication databases.

An e-VRE trust strategy should take into account that there is a high cost associated with the creation and maintenance of extensive metadata and provenance information. The preferred e-VRE strategy is

to collect this metadata automatically as much as possible while allowing users to manually add metadata if they estimate that this is cost-effective. Both the automatic and the manual approach require that the metadata model (CERIF) provides the properties to hold the necessary information.

From a management (rather than technological) perspective, the trust of users in the e-VRE as a whole will benefit from          positive interaction, e.g. through training sessions, newsletters, a transparent and quick response procedure for questions, suggestions and complaints.

## 4.2 Integrity of research (meta)data

It is important to know which versions of both the research data and its metadata were used to produce a result such as a table, chart or publication.  It is also important to know that the data has not been modified, either maliciously or accidentally.  A well established solution to this problem is to compute a cryptographic hash from the data, for example from the SHA family. Cryptographic hashes are at the basis of many modern tools and techniques such at the GIT[20] version management system, Docker[21] application containers, nanopublications,[22] etc. In these systems a particular version may be addressed using its hash key.  Many of the systems allow for associating a hash-key with a name (tag) to facilitate exchange. The use of cryptographic hashes serves two purposes. First of all it allows to unambiguously point at a data collection at a specific version and second it allows anyone with access to the data and the hash to verify that the data is indeed exactly the data that was referred to.

A good application of content hashes is to provide *permalinks* to results. In this scenario the RI creates a document that states how the result was obtained in a machine readable format, together with references (hashes) to the involved data. This document is saved and a hash is created to form the basis of the permalink.  This technique combines well with GIT as well as with nanopublications.  The PROV vocabulary[23] can be used to describe the provenance of the result.

Content-based hashes work well if all data is accessible as files. This is not always the case.  Sometimes the data is only available as a service. In this case we depend on the service to provide us with a reference to a version. This is not always possible; consider for example research results based on the *normalized Google distance[24],* a metric that depends on non-reproducible Google search results. As an alternative we recommend to provide a reference to the service together with the date it was accessed, making clear that rerunning of the service at a later time may produce different outcomes.

# 5 Strategy with respect to legal issues

Any strategy with respect to legal issues—most prominently licensing and intellectual property rights (IPR)—must be fully compatible with the policies of the underlying e-RIs. Broadly speaking, all e-infrastructures in Europe seek to uphold the same basic principles of Open Access for data and services and Open Source for tools and resources. There is also however an emphasis on attribution; that individuals and organisations responsible for procuring data made available via the research infrastructure are properly credited. Beyond that the extent to which e-RIs have formally defined their operational policies regarding IPR and data licensing depends on their implementation state—e-RIs in their preparatory phases are unlikely to have fully formed data access policies. As such, we draw much of our analysis of best practice from those e-RIs that have formed ERICs (European Research

---

[20] https://git-scm.com/

[21] https://www.docker.com/

[22] http://nanopub.org/

[23] https://www.w3.org/ns/prov

[24] https://en.wikipedia.org/wiki/Normalized_Google_distance

Infrastructure Consortiums), which are distinct legal entities and therefore have had need to produce formal statutes and/or access policies. The primary sources of information used in this section were[25]:

- "ICOS data policy, ISIC approved version" *(May 2013)* [26].
- "Statutes of the e-science and technology European consortium for biodiversity and ecosystem research (LifeWatch ERIC)" *(July 2013)*.
- "Statutes of the European Multidisciplinary Seafloor and water column Observatory- European Research Infrastructure Consortium (EMSO ERIC)" *(December 2013)* [27].
- "Commission Implementing Decision of 5 May 2014 on setting up Euro-Argo Research Infrastructure as a European Research Infrastructure Consortium (Euro-Argo ERIC)" [28].
- "EPOS data policy and access rules" *(June 2014)*.
- "ENVRIplus initial data management plan" *(May 2015)* [29].

In brief, the key requirement for an e-VRE is to properly account for the provenance of e-RI resources in order to:

1. Correctly attribute organisations and individuals for their contribution to research outputs.
2. Ensure that any and all terms and conditions of the e-RI are met during operation through the e-VRE.

In principle, as long as there is adequate accounting and correct handling of data, models and services in accordance with their associated licences and the policies of the underlying e-RIs, then most legal issues arising from the use of an e-VRE can be 'pushed back' to the e-RI serving the resources.

## 5.1 Intellectual property rights

Intellectual property rights in this context principally concern three factors:

1. The protection of data products, tools and services developed by the e-RIs to help with the realisation of their services.
2. The protection of resources provided through the e-RIs where those resources originate from the organisations, data centres and individuals that contribute to the e-RI.
3. The protection of the publications produced by users of e-RI services.

An e-VRE must conform to the IPR policies of the e-RIs that it provides a service layer for. Thankfully, it seems that most mature e-RIs in Europe (exemplified by those which have ERICs or are in the process of acquiring them) are converging on very similar IPR and data access policies. From analysis of relevant documents (i.e. the primary sources of information mentioned above), a number of generic observations can be made:

- A number of e-RIs (e.g. EMSO and ICOS) explicitly define intellectual property according to Article 2 of the Convention Establishing the World Intellectual Property organisation (Stockholm, July 1967)—this can probably be taken as the standard definition of intellectual property recognised by the e-RIs.
- The vast majority of data is open for use by almost anybody (subject to due attribution), and primary data (i.e. raw data and measurements) are not protected by copyright, nor are the

---

[25] The statutes of the LifeWatch ERIC and the EPOS data policy were acquired via personal correspondence with members of the respective e-RIs.

[26] http://www.socat.info/upload/ICOS_data_policy.pdf, retrieved 27th May 27, 2016.

[27] http://www.emso-eu.org/site/archive/EMSO-ERIC-statutes.pdf, retrieved 27th May 2016.

[28] http://www.euro-argo.eu/About-us/The-Research-Infrastructure/Statutes, retrieved 27th May 2016.

[29] http://www.envriplus.eu/deliverables/, retrieved 27th May 2016.

ideas and principles underpinning programs in software or hardware (including algorithms and data structures), though the actual code itself is (European Parliament directive 2009/24/EC). Investment in the construction of catalogues and databases is somewhat protected via some *sui generis* rights conferred by the Database Directive (European Parliament directive 96/9/EC), though the scope of protection given to e-RIs is unclear due to a lack of harmonisation across the EU on the subject.

- Most European e-RIs, as might be expected, adhere to European Union policies regarding data access and IPR, including those defined by the INSPIRE directive (2007/2/EC); for example, the LifeWatch ERIC shall "follow European Union policies on data access and IPRs developed under the European Commission Recommendation on access to and preservation of scientific information of 17.07.2012 and its amendments and related instruments" (LifeWatch 2013). It can be assumed that this applies to most if not all ERICs; "with respect to questions of Intellectual Property, the relations between the Members will be governed by the national legislation of the Members and by international agreements to which the Members are parties" (Euro-Argo 2014).

- In general, e-RIs assert ownership of all intellectual property rights created or obtained in the course of their activities via their respective ERICs. Some ERICs, such as Euro-Argo, explicitly note however that intellectual property rights generated by members (meaning member countries) or observers are retained by those entities. Similarly, ICOS notes that intellectual property generated by contributing networks not part of the ICOS Carbon Portal (the central portal by which to access ICOS services) also remain the property of those networks, and some contributors to the ICOS e-RI may also retain IP control subject to the contracts with the ICOS ERIC. In the case of ICOS, it is also noted that third-party IP rights are not automatically accessible to the ICOS e-RI or ERIC, and that ICOS respects the IPR of external modelling groups that have made their derived data products available via the Carbon Portal.

From the above, we derive that an e-VRE should **pass on terms and conditions** required by e-RIs of its users to the e-VRE user, **pass on licensing information** to the e-VRE user so that they are made explicitly aware of the constraints attached to the research assets they are employing, **pass on security credentials** from the e-VRE user to the e-RI where required, and **ensure that its own operation does not violate usage restrictions or expose sensitive information to any party** (e.g. e-RI, e-VRE user, or other third party) **not permitted access**. Moreover, an e-VRE should acknowledge the rights of not only the underlying e-RI, but also the constituent data networks providing resources to the e-RI.

## 5.2 Software and data licensing

Due to a general reliance on open source technologies and open access principles, software and data licensing from e-RIs imposes few constraints on e-VREs or e-VRE users; there are however certain important restrictions regarding attribution, and also in some cases commercialisation.

Generally e-RIs require that users agree to terms and conditions governing access. These tend to be fairly unrestricted, but may prohibit commercial use in some cases, or require that researchers access the data from an IP address mapped to a physical location within a nation contributing funds to a given e-RI; in the latter case, an e-VRE would probably need to share the originating IP of the user interacting with the e-RI's resources. It may also be necessary for e-VREs to pass on terms and conditions from an e-RI service to the e-VRE user before the service can be exploited, and it may be necessary for a given user to actually register with the e-RI, or possibly to have an VRE act as a proxy registrant on the e-RI (which may have limited capabilities compared to a direct registration). Ideally, it should be possible to access and browse the data usage policies for every dataset accessed by the e-VRE.

A number of generic observations about how e-RIs consider data licensing:

- A common licence used for data products and tools developed using scientific data collected by e-RIs appears to be the **Creative Commons CC BY licence**[30], which dictates that attribution is required, but which otherwise permits free use of the data without restriction or discrimination. Other variants (CC:BY:SA requiring any derivative works to maintain the same licensing restrictions, CC:BY:ND, preventing the modification of products/tools, and CC BY-NC, preventing commercial use) may also see use in some cases. Some further licence customisation may occur as e-RI services evolve, leading to the production of e-RI-specific data licences, however it is expected that the same core principles will remain.
- Part of the purpose of ERICs is to "protect Data Providers/author's' … right to the proper acknowledgement and citation, and relieve Data Providers/authors from any legal responsibilities on their behalf" (ICOS 2013); any e-VRE should extend this protection.
- While all e-RIs assert a preference (and intention) for free access to data, *libre* and *gratis*, most if not all ERICs retain for themselves the right to charge for certain aspects of their service, albeit usually restricted to covering operational costs only, presumably as a long-term sustainability contingency measure.
- There is some provisioning for cases when resource providers would withdraw from e-RI consortia. For example, "if, for any reason, a Data Provider (ICOS National Networks and the ICOS CFs) is withdrawing from ICOS e-RI, the Data Provider has a responsibility to give to ICOS ERIC a free of charge, perpetual, non-exclusive, non-transferable right to use the ICOS Data Related Tools (i.e. are the codes, algorithms and software used to generate, collect and process ICOS Data) and all necessary documentation to use the tools in order to archive and process ICOS Data to meet the ICOS e-RI objectives" (ICOS 2013). In principle, such policies can permit continuity of service at the e-VRE level if resources or assets are transferred from one part of an e-RI to another. Nevertheless, it may be necessary to change the policy associated with some resources, and this needs to be pushed to users of those resources via the e-VRE.
- Different conditions may be needed for accessing data from different e-RIs or even a single e-RI. For example ICOS makes a distinction between its own core ICOS data and 'background' and 'sideground' data produced prior to or in tandem with ICOS operations that may be handled or archived by parts of the ICOS e-RI, but which are not bound to the same conditions as ICOS data unless formally integrated into the ICOS data portfolio. This notion of 'auxiliary' data that might be accessible via an e-RI, but is not part of the 'core' contribution may be a recurring one for distributed e-RIs involving facilities and organisations that contribute to multiple different e-RIs for slightly different purposes. As a consequence, it will be prudent to avoid a "one size fits all" policy for e-RI data accessed by a single e-VRE, even for data originating from a single e-RI.

At present, the licensing of software via an e-RI seems to be a mostly unexplored area. While most software produced within e-RIs is considered part of the e-RI's intellectual property, this software is often considered to be part of the e-RI's services, rather than a resource provided to researchers, and so the issues surrounding the licensing of software are not as fully considered as for those of data.

## 5.3 Accounting in e-infrastructures

The implementation of facilities for the accounting of distributed hardware infrastructures may be a necessary condition of close integration between an e-VRE and the e-RIs providing resources to it.

Accounting for infrastructure requires recording every access and interaction with e-RI resources. Where possible the provenance information provided by the e-RIs should be retained, otherwise the e-VRE should at least track what information it can derive (e.g. source URI for data accesses). Most e-RIs disseminate data either by making it available at their own data centres, or by pushing them out to

---

[30] https://creativecommons.org/licenses/by/4.0/

certain international bodies and data networks. In the latter case, there is usually a protocol for tracking the original provenance of data; any e-VRE interacting with international networks should extract that provenance if possible, though it may not always be essential if the network itself is credited. The use of persistent identifiers assists in attribution and accounting by essentially providing a global access point for basic provenance about a research asset. One difficulty lies with the differing aims of accounting and provenance; while a unified model is desirable, some provenance services may not be wholly suitable for accounting, and in any case will likely extend in functionality beyond the needs of accounting.

The e-VRE must provide a transparent mechanism by which to reveal the user driving the data in all cases where registration of users is required. EPOS identifies three classes of users in its "data policies and access rules" (EPOS data policy and access rules 2014):

- *Anonymous* users without identification or accreditation. Only fundamental services should be provided by an e-VRE without some form of registration interaction with the underlying e-RI.
- *Registered* users that have identified themselves via a prior registration process. Such users would expect greater access to e-RI services via the e-VRE, but may not be authorised for complete access. Ideally the e-VRE should provide a mechanism for handling registration, but otherwise there needs to be a way to transfer registration information to the e-VRE so that properly accounting can be carried out.
- *Authorised* users that have not only registered, but have been conferred special privileges allowing access to a greater range of resources. Authorisation can of course be partial given the range of possible privileged resources provided by a number of e-RIs. Again, there needs to be a mechanism by which the e-VRE can provide authorised user credentials to the underlying e-VRE.

In connection to the above, the EPOS data policy and access rules (2014) also identifies three classes of data, which can be considered generally applicable to many e-RIs that deal with a variety of data types:

- *Open* data freely available to users for either direct use or download.
- *Restricted* data available subject to certain conditions; such restrictions could also include the requirement that fees must be paid.
- *Embargoed* data available to certain authorised users for a certain period, before becoming open or restricted.

Generally speaking, the access to restricted or embargoed data will be contingent on the provision of suitable user credentials. The general policy of openness suggests that knowledge of the existence of restricted/embargoed data is not, by default, itself restricted, but if this proves to be the case (most likely for embargoed data), then there needs to be some filtering or partitioning of resource catalogues based on the privileges conferred to a given user. Access to software or resources can be restricted or embargoed using a similar model.

## 5.4 Liability

Generally, e-RIs do not accept liability for any harm caused by the use of their services. EPOS specifically has asserted the need for professional indemnity insurance for EPOS-ERIC employees in their work. In the context of e-VREs, the notion that the source e-RI is not responsible for events arising from the use of their data and tools is well understood by most researchers, and should be made clear in their terms and conditions of use in any case. This should be explicit in an e-VRE security policy.

The e-VRE should avoid shifting any liability onto the user by consequence of failure of its security or ability to correctly restrict access to privileged data, though if access to an e-RI is provided via the correct channels, the e-RI's own authentication and authorisation services should be sufficient.

# 6 Privacy strategy

The e-VRE should guarantee the protection of both personal research data that is accessed via the e-VRE and personal data about the users of the e-VRE and their actions on the system. Regarding personal data in research data, the privacy policy of the e-VRE should conform to the EU data protection rules (General Data Protection Regulation, Regulation (EU) 2016/679[31]). This regulation will be enforced from 25 May 2018 onwards, and includes the following key points[32]:

a) Easy access: data subjects (for example users) are guaranteed to have free and easy access to their personal data and get understandable information about how their data is being processed.
b) Consent: data subjects will be asked for their consent explicitly.
c) The right to be forgotten: data subjects have the right to request erasure of personal data.
d) Data portability: data subjects have the right to transfer their personal data between service providers.
e) Breach: in case of a data breach, organisations are required to notify both individuals and the relevant data protection authority.
f) Responsibility and accountability: data protection must be designed into the business processes for products and services, and privacy settings are set at a high level by default.

As a consequence, data that is being shared via an e-RI (among a specific research community) may not be suitable to be shared through an e-VRE, if it can then be accessed by a larger community,  for a broader set of purposes, and where it can be combined with a larger set of other datasets. The fact that the e-VRE bridges across several e-RIs poses additional challenges with regard to privacy. In D2.1 on requirements elicitation, it was noted that *"Datasets often require removing privacy sensitive variables from it before publication. [...] Moreover, the combination of data with other sources might still make it possible to track the identity of an individual person, especially when open data are combined with social media data. "* Similar risks emerge from a combination with Open Government Data. The use of Semantic Web standards (which VRE4EIC aims to do for metadata) makes combinations of datasets more likely. The privacy levels of data in an e-RI are not always strict enough for an e-VRE. This results in additional requirements related to resetting access control settings (e.g. to disallow combination of data when an e-RI becomes part of the e-VRE), creating awareness with data providers (that their previous privacy policy might no longer be enough).

The ultimate decision of whether data can be shared through the e-VRE and with what AAAI settings is made by the user that owns the data. The e-VRE should notify data owners of potential privacy issues and provide them with information that can be used to decide about the required measures to protect personal information in a dataset (a description of the e-VRE privacy policy, links to relevant documentation including the European Data Protection Directive, training sessions, etc.). In addition, the e-VRE should provide users (data-owners) with information about AAAI options, to enable the user to restrict access to the data to a selection of users. The general aim of the e-VRE is to adhere to the strict privacy policies in the European Data Protection Directive (thus necessarily limiting the amount of data that can be shared through the e-VRE) while at the same time providing reliable AAAI settings

---

[31] REGULATION (EU) 2016/679 OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation). https://www.eugdpr.org/, accessed January 30, 2018.

[32] http://ec.europa.eu/justice/data-protection/document/factsheets_2016/factsheet_dp_reform_citizens_rights_2016_en.pdf, accessed January 30, 2016.

that can be set per dataset by end-users (thus enabling sharing of personal research data where possible).

The e-VRE itself will also collect and store personal data: usage data of who did what on the e-VRE infrastructure (transaction logs) and personal (CERIF) data such as names, affiliations, project membership, authorship, etc. Also with respect to this data, the e-VRE privacy policy is informed by the European Data Protection Directive. As discussed in Section 3, (1) a user's identity is never made public unless she voluntarily agrees to this; (2) the e-VRE may use incentives to convince people to open up their identities to other users. The incentives could, for example, include more visibility of their work and better traceability of their resources, leading to increased trust of other users in their resources, and therefore to more reuse and ultimately to a higher citation count.

Next to technological solutions, the privacy policy of an e-VRE should contain a terms of use document, explicating how and for what purpose the e-VRE collects and treats personal information. In addition, a protocol is necessary regarding the actions to be taken in the event of a security breach.

# 7 Security strategy

In order to provide a meaningful security strategy, we must take into account the nature of an e-VRE and its relationship with the e-RIs upon which it relies and to which it provides homogeneous access. As already stated in section 1.4, the e-VRE should facilitate access to the data and services of the e-RIs, but cannot provide security where an e-RI fails to do this.

In this sense, while requirements from the e-RIs and from the communities are important to define the e-VRE architecture and functions, there are also a set of "requirements for the e-RIs" that should (1) guarantee that e-RIs can be interoperable with the e-VRE security modules and (2) define the balance of responsibilities and competences between e-VRE and e-RIs that will guarantee the expected level of security of the whole ecosystem (e-VRE + e-RIs).

In this section we will therefore discuss the main model to guarantee the overall security and also lay the basis for the definition of the requirements for the e-RIs in terms of security.

Importantly, such analysis and discussion is the result of constant interaction with and participation in other initiatives, some of which are funded by the EU, to discuss specific aspects of the authentication at European Level (e.g., the AARC[33] project) and others which are now discussing such topics (the already mentioned ENVRI and EPOS). In this sense, the VRE4EIC project wants to maximise the synergy among such initiatives, optimise resources and build upon the already existing results.

## 7.1 Handling security in a e-VRE

When dealing with the challenging topic of security in a wide and heterogeneous scenario, it is important to understand the actors, their roles and have an architectural overview of how the different actors will interact. In addition, the three dimensions (or key requirements) that we can consider to study, analyse and design the security of an e-VRE system, are authentication, confidentiality, and access control. In the following section we will discuss the scenario, the requirements to the e-RIs and candidate technical solutions.

## 7.2 Architecture of the e-VRE / e-RI security system

Authentication enables the system to know who is using the system and on the basis of this knowledge the system can provide proper permissions. From another point of view, confidentiality should guarantee that provided data or services are not accessed by users who do not have the proper rights. Access control allows restricting access only to users who have the rights to do so, and may be used to apply special conditions to the provided resources.

A number of questions may help to introduce the first and main topic:

- Should the e-VRE build a new security system?
- How to manage existing security system within the e-RIs?
- How should the e-RIs' security systems interoperate with the e-VRE security system?

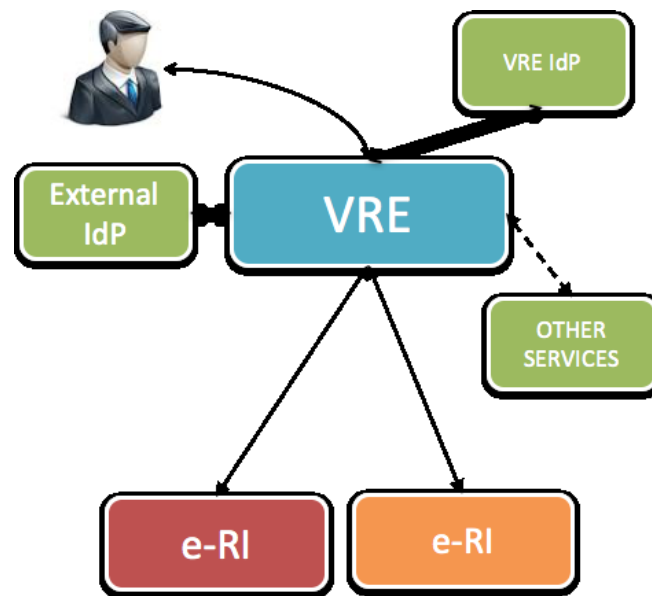Answering such questions means understanding the overall architecture of the e-VRE / e-RIs.

---

[33] https://aarc-project.eu/

Figure 1: access to a e-VRE and underlying e-RIs[34].

In the scenario presented in Figure 1, the e-VRE provides a single access point so that the user has a homogeneous view of heterogeneous resources. This single access point includes an Authentication Authorization Accounting Infrastructure (AAAI), by means of which the user must be able to access, with one single authentication, all e-RI resources "brokered" by the e-VRE. In order to simplify the Identity Management, the user will be able to use his/her own credentials, when they are provided by some well known and trusted Identity Provider (e.g. eduGAIN[35]). Once the user is authorised, s/he must be able to access e-RIs resources.

Such scenario implies that the system must:

- Integrate credentials from several Identity Providers (IDPs) (AAAI integration)
- Be delegated to access to e-RIs on behalf of the user (delegation)
- Provide IDP functionalities, when the user has no pre-existing credentials


It implies the following requirements to the e-RIs:

1. The e-RIs should provide an AAAI system.
2. The e-RIs AAAI systems should enable the e-VRE to securely access their resources.
3. The security of the resources provided by each e-RI must be handled at the e-RI level. It means, for instance, that data storage and backup is handled by e-RIs, and that access to them should be regulated by the e-RI. Of course both e-VRE and e-RIs should be in a trusted domain where the e-RIs trust all users / connections from the e-VRE. In addition, secure data transmission (e.g. through encryption) is a shared responsibility of the e-RIs and the e-VRE.
4. The AAAI systems provided by e-RIs should be standard and interoperable with the e-VRE.
5. The different types of authorization should be somehow homogenised at e-RI level. It means that a limited number of authorization "groups" at e-VRE level should fit to all e-RIs.

---

[34] In the framework of the EGI-ENGAGE EPOS competence Center (https://wiki.egi.eu/wiki/EGI-Engage:Competence_centre_EPOS), such architecture is being tested. More details at https://wiki.egi.eu/wiki/File:EPOS_Competence_Center_USE_CASE-DS_.doc

[35] http://services.geant.net/edugain/

## 7.3 Existing technical solutions from main e-RIs / use cases

The e-RI characterizations evidenced that one of the driving use cases - EPOS - is converging toward a technical solution that tackles some of the main issues of the AAAI.

Such a solution is supported by several initiatives, namely EGI[36], EUDAT, AARC. AARC is at a different level from the two others though, as AARC wishes to provide an extension to IDPs such as EduGAIN, supporting user attributes and X.509 certificate based services. AAAI solutions basically enable a user to have single authentication to access all resources, and under specific circumstances (depending on the protocol) enable an e-VRE to be delegated to act on behalf of the user .

Such a solution is the UNITY[37] software, which facilitates the establishment of a solution for identity, federation and inter-federation management. Or, looking from a different perspective, it is an extremely flexible authentication service. UNITY is a service that enables login to a web service using various protocols. It supports the LDAP[38] protocol (e.g. OpenLDAP[39] or Active Directory[40]) and authentication can be performed with various identity providers, amongst others the EduGAIN federation previously mentioned. UNITY is open source software licensed under the BSD licence[41].

In order to connect existing AAAI approaches from e-RIs into one e-VRE ecosystem, an AAAI hub is needed which will assure interoperability between existing technologies. HUB technologies passing logins, passwords, and such are deprecated and should no longer be used. Instead attribute based solutions are advised. An IDP should return a digitally signed document that states the identity of the user. Once a user is authenticated within the infrastructure all the authorisations can be done using the attributes only. The term attribute is used here to describe properties of the user, e.g. his/her name, email, affinity, role. A set of the e-VRE specific attributes will have to be defined on the hub.

Current protocols include: LDAP, OpenID, X.509 certificates, EduGAIN (different from the others, it is more a political project using SAML as its core technology).

In the EPOS AAAI Hub, and from the technological point of view, UNITY provides plug-ins for many IDPs (not just the four mentioned above). For instance, a UNITY instance registered in EduGAIN is sufficient to enable EduGAIN authentication from any service attached to UNITY. Attribute management is flexible from the administrator point of view. REST APIs are also available.

A solution like this, while far from being the panacea, can indeed help to tackle at least the main problem, that is to say the federated identity management: in practice it will enable any user (with credentials supported by UNITY) to access the e-VRE.

What still remains an open challenge is the delegation. This topic has been discussed in EPOS, EGI, EUDAT, AARC, and still remains a work in progress.

In the framework of our project, a draft study and architecture may provide important input to other initiatives working with the specific AAAI technologies. They could indeed align their future developments based on e-VRE requirements, moving towards the integration, interoperability and optimization of resources.

## 7.4 Two-factor authentication

Two-factor authentication (2FA) is a method of confirming a user's claimed identity by utilizing a

---

combination of two different components (factors). Two-factor authentication is a type of multi-factor authentication. Using 2FA, users can authenticate using something that only the individual user knows (for instance login and password) plus a one-time-valid, dynamic passcode, typically consisting of a number of 4 to 6 digits. The use of two-factor authentication is strongly recommended for an e-VRE if privacy-sensitive data or metadata are being accessed.

When choosing a 2FA method, one should take into account security (e.g. whether messages can be encrypted), privacy (e.g. whether additional user information needs to be stored) and ease-of-use (e.g. the number of (additional) devices that a user must carry). The typical implementation of 2FA uses mobile phones as recipients of the second factor: the code is sent to the mobile phone of a user via SMS or push notification. The advantage of using a mobile phone as the second factor is that there is no need to provide to a user a dedicated token that must be used to complete authentication.

For the e-VRE prototype system that is currently being developed within WP3, we have chosen Telegram as a local 2FA solution. Telegram[42] is a messaging app available for free and released for most common mobile operating systems (Android, IOS, WP). It supports encrypted communications and provides APIs that can be used by third-party developers to integrate Telegram functionalities with their applications. An advantage for the privacy of the user is that Telegram doesn't need the phone number of the user. Additionally, since Telegram clients are available for many kind of devices (smartphones, PCs, tablets), the second factor is not tied to a single specific device: a user can install the Telegram app to multiple devices and synchronise all installations, the passcode will then be received in every device having the Telegram client installed. It is important to note that the e-VRE Telegram Bot is not only used for the 2FA, it has been designed to be extended in order to enable users to access a number of e-VRE functionalities (details are described in Deliverable 3.3).

The e-VRE prototype system interacts with the Telegram Bot API[43] to implement the 2FA. Essentially, the second factor of the 2FA is managed by a *Telegram Bot*. Telegram Bots are special applications that can be connected to Telegram accounts and that do not require a phone number to be used. A connection between a user and a Telegram Bot is uniquely identified, and the identifier will persist until the user removes the Bot from the list of Telegram contacts.
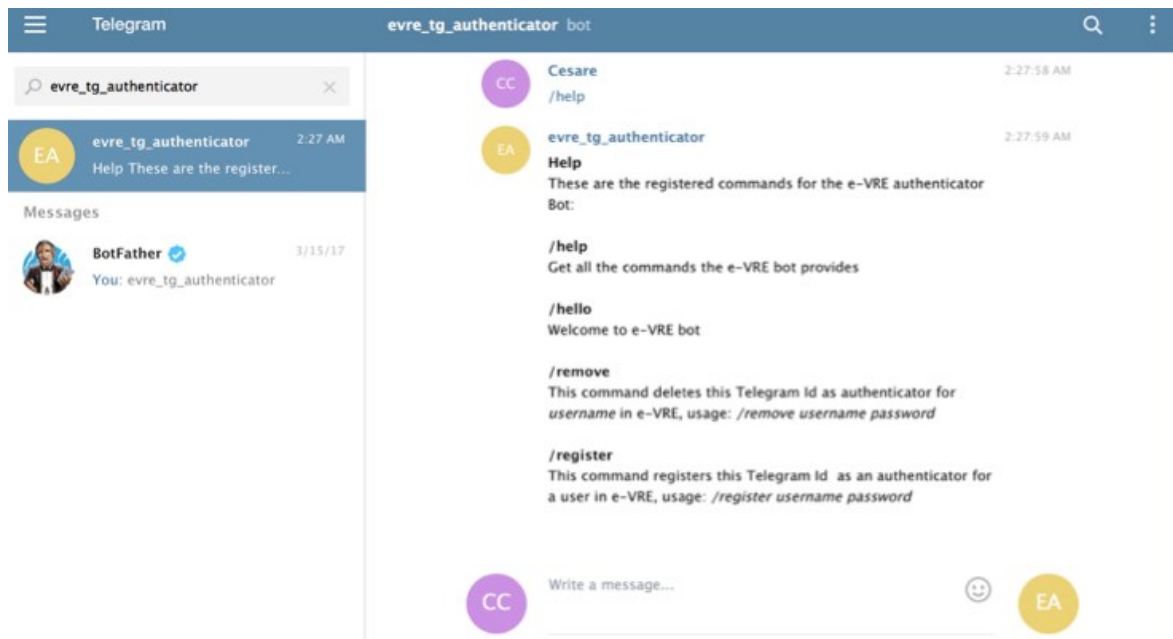
A user that wants to use the Telegram channel to exchange passcodes with the e-VRE, can add the account of the e-VRE authenticator Bot (called *evre-tg_authenticator*) to the Telegram contact list. The screenshot below shows how the VRE4EIC Bot may appear on the Telegram app installed on the smartphone of a user. When an e-VRE user sends the **/register *username password*** command to evre-tg_authenticator Bot, the Bot gets the identifier of the channel and stores it in the user profile, enabling the channel for the exchange of the passcode for 2FA. The command **/remove *username password*** disables the 2FA channel. The authentication of a user that has registered its username on e-VRE via evre-tg_authenticator occurs as follows:

- The user, by interacting with a Web UI, starts the 2FA protocol sending the login/password credentials to the e-VRE.

- The e-VRE first checks the validity of the credentials, then establishes the id of the Telegram channel that the user wants to use to receive the second factor.

- The user receives a first answer, again a web UI, informing her/him that a passcode is being sent via the registered channel.

- The e-VRE system generates a one-time-valid passcode and sends it to the requesting user, via the Telegram channel.

- The user receives the second factor (the passcode) in the Telegram app and uses the web UI to send it back it to the e-VRE system.
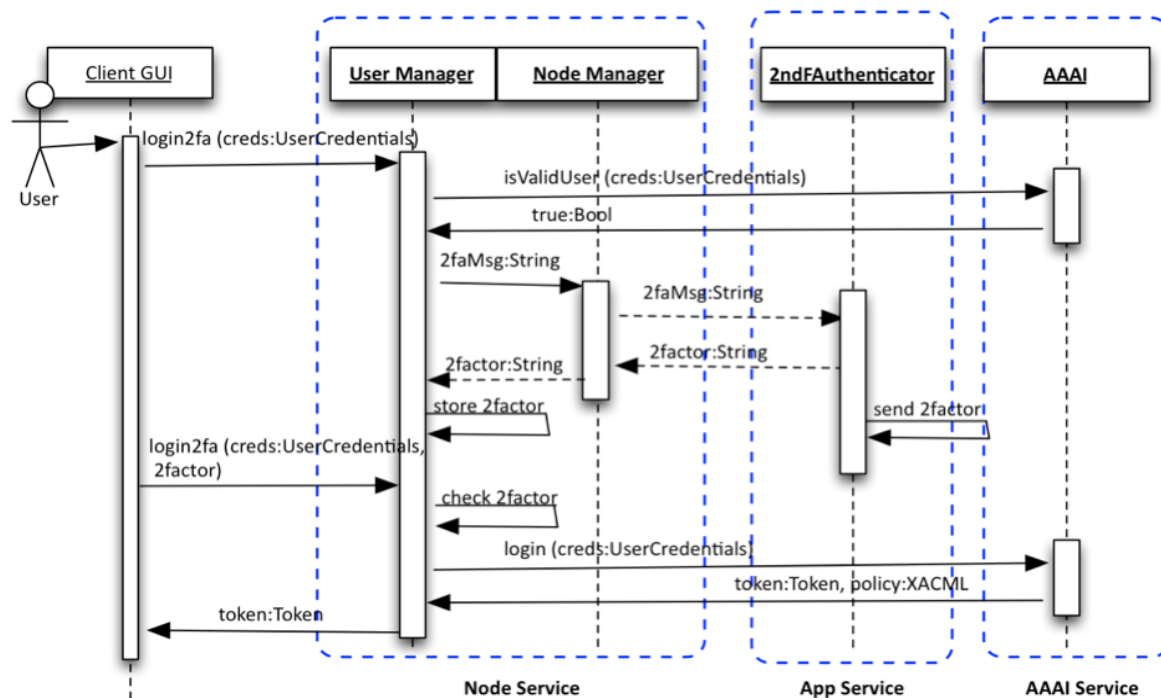
---

[42] https://telegram.org/

[43] https://core.telegram.org/bots

- The credentials previously sent and the passcode are used by e-VRE to actually authenticate the user.



**Screenshot of the Telegram messaging app, showing the VRE4EIC bot for authentication.**

From an architectural point of view, the 2FA protocol is implemented by the cooperation of a number of components. Figure below shows a UML sequence diagram that describes how these e-VRE components interact to implement it. A detailed description of the diagram is presented in Deliverable 3.3. As shown in the diagram, there is a specific component for the management of the second factor. It interacts asynchronously with e-VRE main components. This architectural approach enables us to easily replace/update this component and also to have several different channels.



**UML sequence diagram with interaction between the 2FA and AAAI components.**

# 8 Metadata Strategy

The metadata strategy will depend on the CERIF metadata model. In CERIF, a user is a person represented by the CERIF entity "cfPerson", that may have values for properties such as Birth date, Gender, URI, First / Family / Other names, Research interest, and Keywords.
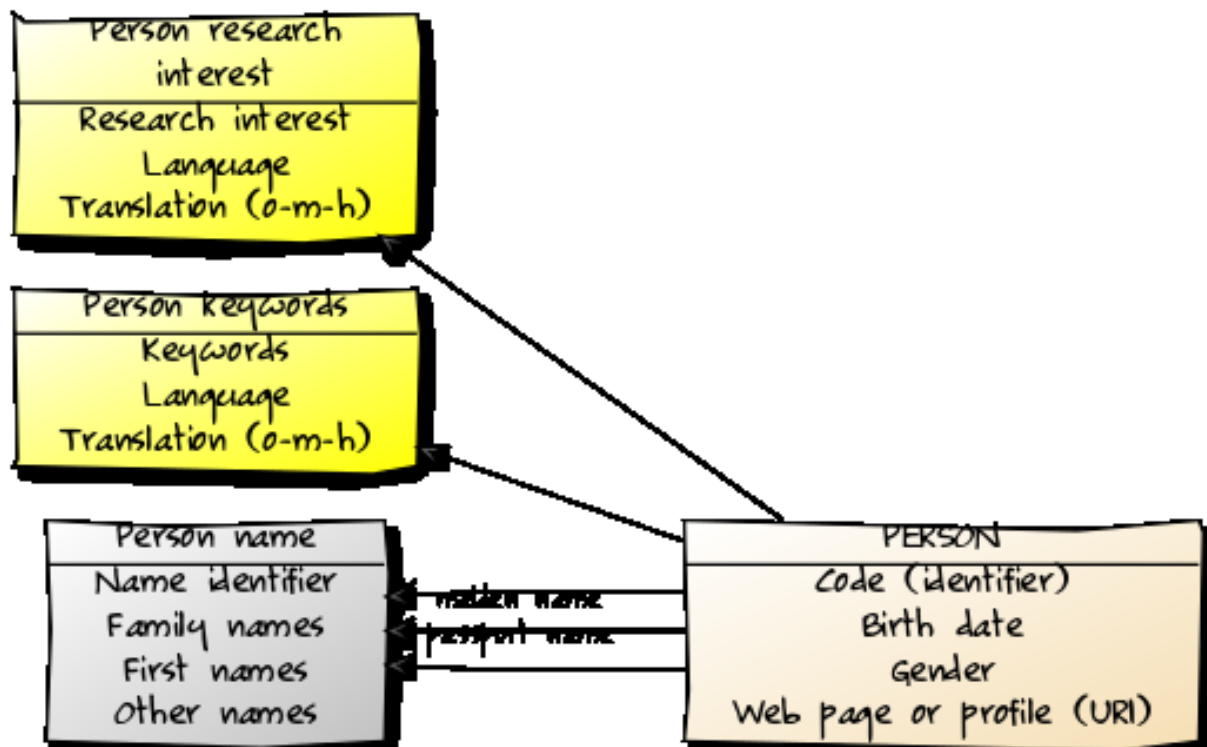


**Figure 2: Representation of a person in CERIF, used to model a User.**

## 8.1 Role between User and (each) Resource

One or many Role(s) can be given to the User / Person as a direct classification of this Person, and/or to the relation(s) between this Person and the resource to be accessed. In Figure 3 below, the Roles are defined by the pair of properties marked with * (cfClassId, cfClassSchemeId).

The role can then define the access level of the user to the resource (e.g. Creator, Editor, Viewer). That would be implemented through a vocabulary (cfClassScheme) and terms (cfClass) defining these levels of access. To use this approach, we need to create and store the relations between all users and all resources, which may not be a scalable solution. Also, all these relations need to have a role, which cannot be managed humanly, so we would need to rely on a default value (for example, a given resource can be viewed by all users, so the role of Viewer is set; for other resources, the default value could be "No access"). Then, only specific and less numerous values would be changed manually upon need. So, this approach is applicable only if the total amount of relations generated by the numbers of users and resources is technically manageable, and the need for manual definition is humanly manageable.

cfPerson_Equipment

cfPerson_Facility

cfPerson_ResultProduct
(datasets, software components)

cfPerson_ResultPublication (any
document)

cfPerson_Medium (other
media than document)

cfPerson_Service (software
services including workflows)

**Figure 3: CERIF Linked Entities involved in User-Resource relations**

## 8.2 Separation of role and access permission

To avoid the scalability issue and the necessity to manage all users-all resources access definition, we propose to use Role-Based Access Control (RBAC) which separates several layers: users and groups on one hand, and roles, permissions and resources (or types of resources) on the other hand. Then users or security groups can be mapped to roles. The Wikipedia entry on RBAC summarizes the design adequately[44]:
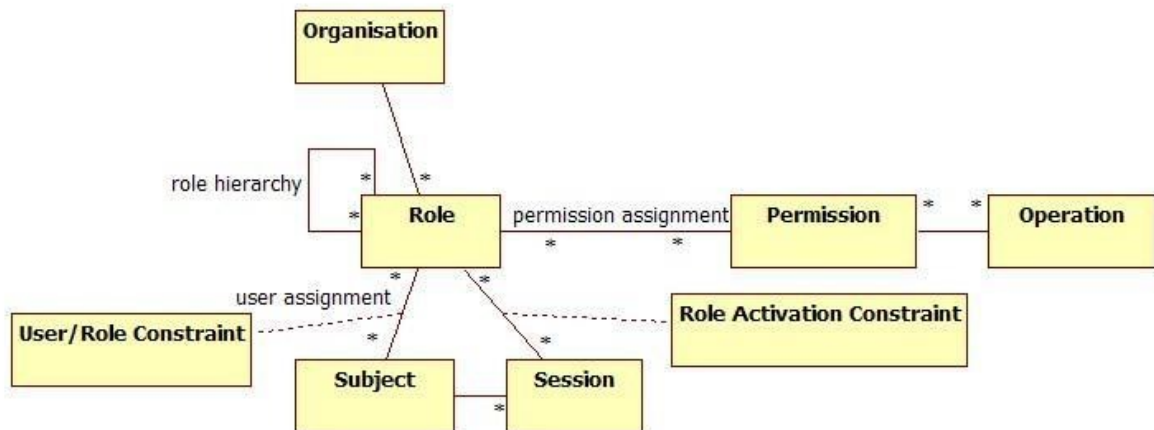
---

**Figure 4: RBAC model[45]**

*"When defining an RBAC model, the following conventions are useful:*

- *S = Subject = A person or automated agent*
- *R = Role = Job function or title which defines an authority level*
- *P = Permissions = Approval of a mode of access to a resource*
- *SE = Session = A mapping involving S, R and/or P*
- *SA = Subject Assignment*
- *PA = Permission Assignment*
- *RH = Partially ordered Role Hierarchy. RH can also be written: ≥ (The notation: x ≥ y means that x inherits the permissions of y.)*
- *A subject can have multiple roles.*
- *A role can have multiple subjects.*
- *A role can have many permissions.*
- *A permission can be assigned to many roles.*
- *An operation can be assigned many permissions.*
- *A permission can be assigned to many operations.*

*A constraint places a restrictive rule on the potential inheritance of permissions from opposing roles, thus it can be used to achieve appropriate separation of duties. For example, the same person should not be allowed to both create a login account and to authorize the account creation."*


With the CERIF concepts, that would be represented as follows:

- A classification of users or security groups (as organisation units): "cfClass for role", that can be multi level defining a hierarchy of roles ("Team leader", "Manager", "Director",…);
- A relation between the "cfClass for role" and resources (such as facilities, equipment,…);
- A Permission as the semantic (another cfClass, "View", "Edit", "Delete") of the relation between the Role and the Resource.
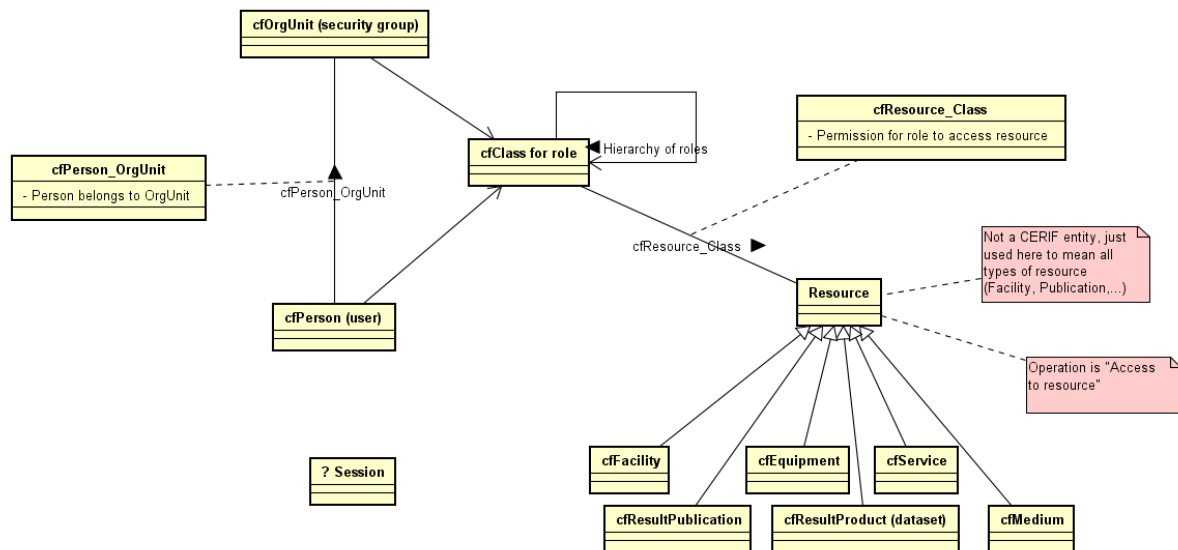
---

**Figure 5: Permission between Roles and Resources**

Again, the need to map all roles to all resources is probably not scalable nor manageable. Also, the "Create" permission cannot be defined. So, the Permission should be between a Role and a Resource type.
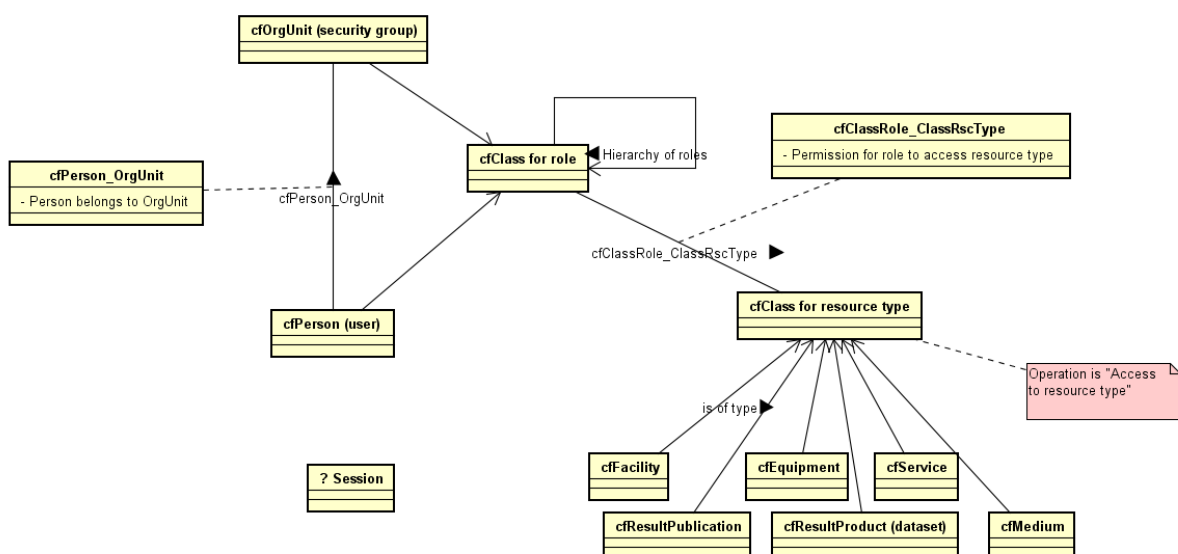


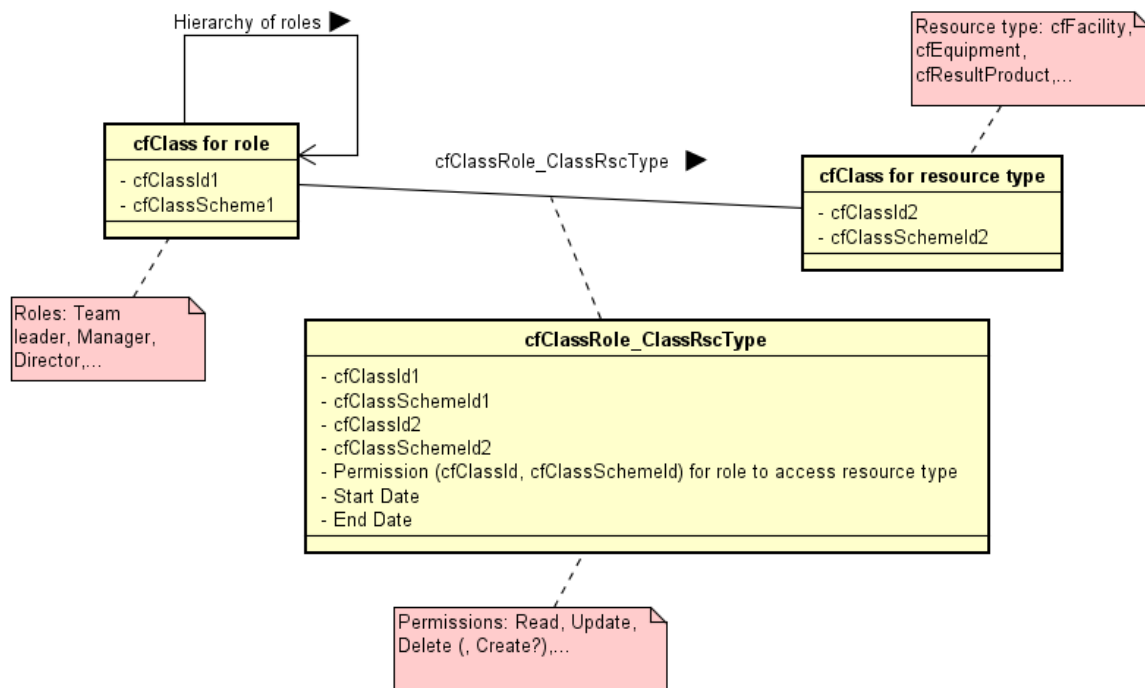**Figure 6: Permission between Roles and Resource types**

**Figure 7: Representing the Roles, Resource types and Permissions as CERIF cfClass and cfClass_Class.**

## 8.3 Persons with different users (and roles) in different contexts

There are other models of role-based access to cover specific needs: OrBAC for organisation-Based Access Control, TMAC for Team-based Access Control (for applying role-based access controls in collaborative environments), TBAC for Task-based Authorization Controls[46]. In addition, there may be situations where a person changes roles, affecting the needed update frequency of the metadata. If the need arises, the CERIF metadata model can be extended to cover this case.

## 8.4 Quality assurance strategy for AAAI metadata

Metadata as discussed above, as expressed using the CERIF model, will be eventually used to directly make automated authentication and authorization decisions. Errors in the metadata, or outdated metadata, can thus have severe consequences and breach the security of the underlying systems (which will be many, assuming a single sign-on architecture) and compromise the protection of sensitive data stored within these systems. Any strategy on AAAI metadata needs to include explicit procedures to assure the quality of that metadata. Special attention is needed to raise awareness in all parts of the organisation that deal with this metadata of the importance of its quality.

Especially research organisations that are already familiar with the use of CERIF metadata in the context of their Current Research Information System (CRIS) need to be aware that while the CRIS metadata and AAAI metadata may look very similar, their intended use is very different and the impact of mistakes in the data is very different as well. Any AAAI strategy needs thus to take explicit measures to avoid confusion between these datasets. For example, while organisations may choose to automatically populate their CERIF AAAI records with user identity information from their current login metadata, it is unlikely that this data already provide the attributes required to enable role-based

---

[46] http://orbac.org/?page_id=4 HYPERLINK "http://orbac.org/?page_id=4"

access (RBAC). Especially when the organisation's CRIS system does contain this type of information (e.g. which researcher is affiliated with which project), it will be tempting to populate the AAAI records directly from the (potentially self-assigned) CRIS data, with all the security risks involved if the quality assurance of the CRIS data is insufficient.

It is thus of key importance that it is clear who is responsible for the quality assurance of the AAAI metadata. In many organisations, CRIS metadata is often entered directly by the researchers involved, and loosely maintained by the library or communication department for consistency and completeness. In contrast, it is typically the IT department that is responsible for assuring the quality of traditional AAAI metadata (like password files or LDAP records), even if they use data from other departments as input. Our recommendation to organisations that participate in the VRE4EIC federated AAAI infrastructure would be to put the responsibility of AAAI metadata quality with the same department that is responsible for the organisation's own AAAI records, and to make this decision independent of the data model used to represent this data.

## 8.5 Mapping of policies to CERIF

Policy as a document

cfResultPublication

Organisation issuing the policy

cfOrgUnit

Issuer, date of issue

List of

cfClass_Class

Permission

cfClass

Resource type
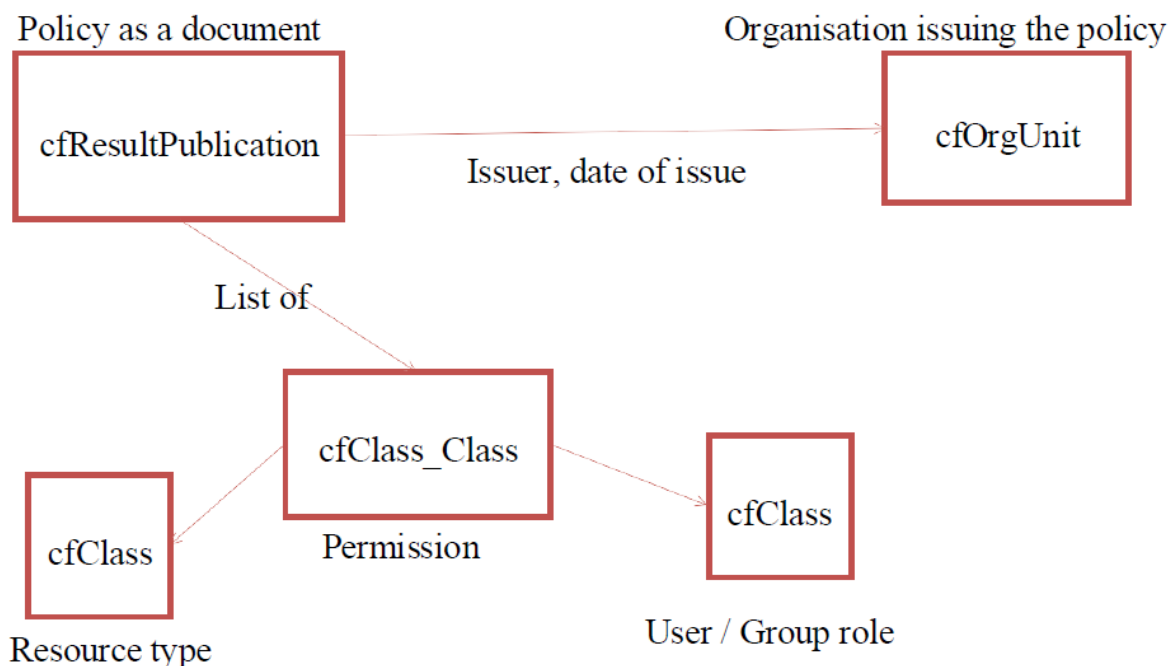
cfClass

User / Group role

Figure 8: making the relation between policies and organisations explicit with CERIF

Where policies come in the form of documentation (as is the case in, for example, a data management plan), it is recommended that the e-VRE uses a standard format to make explicit which organisation is issuing the policy and, where possible, also to model what the policy entails. As a first descriptive approach, a policy is a document linked to an organisation as its Issuer, with a given date of issues, and possibly other dates (of approval, publication,...). A second step includes describing the policy as a list of permissions, using the model described before, where a permission is a classification of a relation between a resource type (itself a classification) and a user or group role (itself a classification). Storing this list of permissions allows automated checking of policy compliance. The above image illustrates how this can be done with CERIF.

# 9 Mapping of strategy to architecture and software components (T3.1 and T3.3)

The e-VRE Reference Architecture has been presented in Deliverable D3.1. This section summarises the fundamental characteristics of the Reference Architecture and describes how the strategies described above are taken into account by the Reference Architecture.

The AAAI (Authentication, Authorization, Accounting Infrastructure) component of the Reference Architecture is designed to cover the strategies presented in the previous section concerning trust, security and privacy functionality required for the proper operation of the e-VRE. To this end, the AAAI component provides four interfaces:

- Authentication interface: providing methods to check credentials in order to authenticate users and agents on the e-VRE.
- Accounting interface: providing methods to track operations and keep information about system states.
- Authorisation interface: providing methods enabling other components to check permissions for executing operations on external resources. As reported in [1], once a user or an agent is authenticated within the e-VRE infrastructure all the authorisation should be done using Attribute-Based Access Control (ABAC) only; operations in this interface are used to manage attributes.
- Cryptographic interface: providing methods to encrypt/decrypt communication with e-RI resources. In the design of the e-VRE architecture, the tradeoff between security and performance will need to be dealt with.

These interfaces are used by other e-VRE components in the Architecture in order to check authorisations. These components are:

- User Manager: uses Authentication interface to verify user credentials, Accounting interface to log access and user management operations, can use Cryptographic interface when encryption is required to interact with agents, uses Authorisation interface for Events management or other tasks.
- Resource Manager: uses Authorisation interface for operating on e-RI resources, can use Cryptographic interface for encrypted interactions.
- Query Manager: uses the Authorisation interface for interactions with e-RI storages and datasets.
- Model Mapper: can use Encryption interface for implementing mappings.
- App Manager: uses Authorisation interface, Accounting interface and possibly Encryption interface to interact with e-RI resources.
- Metadata Manager: can use Encryption interface and Authorisation interface.
- E-VRE Web Services Component: uses Authorisation interface.
- Workflow manager: uses Authorisation interface for verifying permissions to use e-RI resources, Accounting interface to track operations, can use Encryption interface.
- Linked Data Manager: uses Authorisation interface.

The reader is referred to D3.1 for a detailed account of the Reference Architecture components, their interfaces and the signatures of the methods that these interfaces comprise. A major goal is to avoid significant performance issues in processing attributes for authorisation requests and for evaluating related authorization policies.

The microservice-based Technical Architecture that has been devised to implement the Reference Architectures provides a AAAI Service that realizes the four interfaces discussed above. The Technical Architecture is presented in detail in Deliverable D3.3. The AAAI Service has been selected by the Gap Analysis (presented in D3.2) as the highest priority building block of the Reference Architecture. As such, it has been implemented and delivered as part of D3.3. The technology selected for the implementation is UNITY, which at the moment is the open-source solution that best covers the interfaces of the AAAI component.

# 10 List of recommendations

*Privacy recommendations*
- PR1: The e-VRE should have a privacy policy that conforms to the European Data Protection Directive.
- PR2: The e-VRE user should be aware of and agree to the privacy policy of the e-VRE.
- PR3: The e-VRE should guarantee the privacy of both users of the e-VRE (authentication and access logs) and of sensitive research data that is stored through the e-VRE.
- PR4: Privacy recommendations with respect to research data management:
  o The privacy levels of data in an e-RI are not always strict enough for a VRE. This results in additional requirements related to resetting access control settings and creating awareness with data providers.
  o The e-VRE should notify data owners of potential privacy issues arising from e-VRE-wide sharing of resources, and provide them with information that can be used to decide about the required measures to protect personal information.
- PR5: Privacy recommendations with respect to e-VRE usage data
  o The privacy policy of an e-VRE should explicate how and for what purpose the e-VRE collects and treats personal information.
  o In some cases there may be a tension between the need to record provenance of datasets, including information on who did what, and the need to protect the privacy of users, including their identities and access logs. The policy needs to address this.
  o The e-VRE should provide functionality to remove usage data from the e-VRE.

*Trust recommendations*
- TR1: The task of the e-VRE is to provide (CERIF) metadata related to trust. At the e-VRE level, the main requirement is to correctly convey the information that is already present at the e-RI level (incl. data ownership, permanence, licensing and liability) of each dataset.
- TR2: An e-VRE must conform with the IPR policies of the e-RIs that it provides a service layer for.
- TR3: An e-VRE trust policy should take into account that there is a high cost associated with the creation and maintenance of extensive metadata and provenance information. The preferred e-VRE strategy is to collect this metadata automatically as much as possible while allowing users to manually add metadata if they estimate that this is cost-effective.
- TR4: An e-RI should implement methods to ensure that a specific version of a dataset or metadata has not changed, i.e. that the original version can be permanently identified, for example by using hashes. The e-VRE should correctly convey the necessary information to permanently identify a dataset (version) that is already present at the e-RI level to the user.

*Security recommendations*
- SR1: The e-RIs form the baseline for security, privacy and trust for the research data they manage; the e-VRE must guarantee standards that are at least as strong as the e-RI.
- SR2: The e-VRE security policy should make explicit who is liable in case of different types of security breaches. In addition, a protocol is necessary regarding the actions to be taken in the event of a security breach.
- SR3: A successful e-VRE is compatible with a wide variety of identity providers in order to suit the needs of associated e-RIs.
- SR4: The e-VRE should be able to pass on security credentials from the e-VRE users to the e-RI.
- SR5: The e-VRE should ensure that its own operations do not violate usage restrictions of resources of the e-RIs.

- SR6: The e-VRE should be compatible with several external access mechanisms and be able to include new ones when new e-RIs connect to the e-VRE, and allow unrestricted access to open data.
- SR7: The e-VRE provides Role-Based Access Control (RBAC) to separate several layers: users and groups on one hand, and roles, permissions and resources (or types of resources) on the other hand. This needs to be enforced for all interfaces.
- SR8: participating organizations should be required to have clear documentation about who is responsible for the quality of the AAAI metadata that is used for RBAC by the e-VRE. It is recommended that the responsibility lies with the same department that is responsible for the organisation's own AAAI records.
- SR9: Multi-factor authentication is recommended for (privacy-)sensitive data. The method of choice depends on three aspects: security (i.e. encryption), privacy (i.e. storage of additional user information) and ease-of-use (i.e. the need for additional devices).
- SR10: Federated authentication and authorization using research institutes' identity and attribute services at European or world scale are preferred. Until the necessary (authorization) services have become commonly available, the e-VRE should implement a local solution in such a way that e-RIs can easily switch to larger-scale solutions in the future.

Table 9 (below) shows which recommendations impact which components of the Reference Architecture. The recommendations will be used to validate the Reference Architecture from the privacy, trust and security point of view, and, if necessary, to amend the interfaces of the components making them compliant with these recommendations.

| Architecture and software components Recommendation | UI | AAAI components | Cryptographic interface | Interoperability Manager | Metadata Manager | Linked Data Manager | e-VRE Web Services |
|---|---|---|---|---|---|---|---|
| PR2 | x | x | | | | | |
| PR3 | | x | x | | | | |
| PR5 | | | | x | | | |
| TR1 | | | | | x | | |
| TR3 | | | | | x | | |
| TR4 | | | | | x | x | |
| SR1 | | x | | x | | | |
| SR3 | | x | | x | | | |
| SR4 | | x | | x | | | |
| SR5 | | x | | x | | | |
| SR6 | | x | | x | | | |
| SR7 | | x | | | x | x | x |
| SR9 | | x | | | | | |
| SR10 | | x | | | | | |

Table 9: Mapping of recommendation to architecture. Recommendations that are not listed refer to documentation instead of software components.

# 11 Summary and conclusions

Security, privacy and trust at the level of the e-RI and the e-VRE are highly intertwined. The e-RIs form the baseline for security, privacy and trust for the research data they manage; the e-VRE must guarantee standards that are at least as strong as the e-RI's. The added value of the e-VRE is mainly to (1) be a single access point for researchers towards resources in several e-RIs, making AAAI issues a focal point of e-VRE design, and (2) aggregate metadata from the various e-RIs, making interoperability a crucial issue.

*AAAI*

Several AAAI solutions are currently in use at the e-RIs characterized by the project. An e-VRE should support as many of the popular AAAI solutions as possible. EPOS uses UNITY for this purpose: a UNITY instance registered in EduGAIN is sufficient to enable EduGAIN authentication from any service attached to UNITY. This can tackle the issue of federated identity management: in practice it will enable any user with almost any credentials (for example those supported by UNITY) to access the e-VRE. What still remains an open challenge is the delegation. This topic has been discussed in EPOS, EGI, EUDAT, AARC, and still remains a work in progress. As already pointed out, UNITY is also employed as underlying technology for the implementation of the AAAI Service in the Technical Architecture of the VRE4EIC project.

Authentication at its most basic level will enable an e-VRE to confirm that a user is who she says she is. However, on top of that, an e-VRE will often need information about the access rights of the user. This information can be represented in CERIF. We distinguish two scenarios here: (1) the identity provider provides this information to the e-VRE on request. This is the preferred solution. However, some identity providers do not support this functionality, and/or legal issues might prevent this. Alternatively, (2) the e-VRE could rely on agreements with individual e-RIs to provide them with the information. In this case, the CERIF information might be stored locally to the e-VRE. This solution is less ideal since the information is less secure and updates will not be automatically incorporated, harming provenance and trust.

*Interoperability*

Interoperability of metadata plays a role in security, privacy as well as trust. Security credentials need be passed on between the e-VRE and the e-RI to authenticate users, as well as metadata about the level of security of a resource (e.g. open data, restricted access, embargoed data). An e-VRE should pass on licensing information and other terms and conditions required by e-RIs to the e-VRE users, so that they are made explicitly aware of the constraints attached to the research assets they are employing.

Trust in a resource relies on knowledge about the creator and the processes used to create it. The strategy of the e-VRE in this respect will be to provide functionality for users to access the trust-related information of the e-RIs when working with a dataset via the e-VRE. This strategy is preferred over the alternative - to provide functionality to create this kind of metadata within the e-VRE environment.

Resources that are being shared via an e-RI (among a specific research community) may not be suitable to be shared through an e-VRE for privacy reasons. The ultimate decision of whether data can be shared through the e-VRE and with what AAAI settings is made by the user that owns the data. The policy of the e-VRE is to request explicit confirmation of the AAAI settings of resources when an e-RI joins the e-VRE. This includes the provision of information that data owners need to decide about the required measures to protect personal information in a dataset (e.g. a description of the e-VRE privacy policy, links to relevant documentation, including the European Data Protection Directive). In addition,

an e-VRE should provide users (data-owners) with information about AAAI options, to enable a user to restrict access to the data to a selection of users. The general aim of the e-VRE is to adhere to the strict privacy policies in the European Data Protection Directive (thus necessarily limiting the amount of data that can be shared through the e-VRE) while at the same time providing reliable AAAI settings that can be set per dataset by end-users (thus enabling sharing of personal research data where possible).

CERIF provides a metadata model for information about users, groups, roles, permissions, resources and types of resources. It can be used to realize the needed interoperability, by mapping local metadata schemas at the e-RI level to CERIF. To avoid the scalability issue and inability to manage all users-all resources access definition, the idea with Role-Based Access Control (RBAC) is to separate several layers: users and groups on one hand, and roles, permissions and resources (or types of resources) on the other hand. Then users or security groups can be mapped to roles. Note that the storage of CERIF data, whether at the e-RI level or at the e-VRE level, brings additional privacy concerns regarding the personal data it contains. An e-VRE strategy should guarantee the secure storage, transmission and backup of these data.

### *E-VRE / e-RI relationship*

In general, the added value that an e-VRE can bring depends on design decisions at the e-RI level. For example, the more interoperable the metadata formats are, the better the e-VRE can provide a user with consistent trust-related information, allowing, for example, a comparison of datasets. The more e-RIs use common AAAI solutions, and the more information these solutions provide, the more likely it becomes that an e-VRE can become a true single-sign-on environment. On the other hand, in order to kickstart the development of e-VREs it seems preferable to be as inclusive as possible and to limit the hard requirements for the e-RIs that want to join. The strategies described in this document are based on incentives, where the e-RIs are incentivized to implement 'the right' AAAI solutions, and users (mainly data owners) are incentivized to disclose their identities to other users. In addition, the described strategies highlight the need for explicit documentation, for example about liability of the e-VRE.

### *Public availability*

Upon acceptance of this strategy document, it will be made publicly available and especially distributed to EPOS and ENVRIplus and other use case partners. The introduction of an e-VRE can help spreading the use of safe and privacy-aware user authentication to e-RIs because it both simplifies this process from the user perspective and from the RI perspective. It brings SSO (Single Sign-On) to the user, avoiding the need to create an account for every RI.  At the same time it simplifies implementing access control for the RI which now only has to implement authentication with the e-VRE and no longer has to deal with user management and the associated privacy concerns.

# 12 References

De Montjoye, Y.A., Hidalgo, C.A., Verleysen, M. and Blondel, V.D. Unique in the crowd: The privacy bounds of human mobility. Scientific reports, 3, 2013.

D. Artz and Y. Gil.  A survey of trust in computer science and the semantic web. Journal of Semantic Web, 5(2): 58-71, 2007.

J. Golbeck. Trust on the World Wide Web: A Survey. Foundations and Trends in Web Science, 1(2):131-197, 2006.

Ceolin, Davide, Archana Nottamkandath, and Wan Fokkink. "Efficient semi-automated assessment of annotations trustworthiness." Journal of Trust Management 1(1): 1-31, 2014.

Castelfranchi C, Falcone R. Principles of Trust for MAS: Cognitive Anatomy, Social Importance, and Quantification. In Proceedings of the 4th International Conference on Multi-Agent Systems (ICMAS), IEEE Computer Society pp 72–79, 1998.

J. Sabater and C. Sierra. Review on computational trust and reputation models. Artificial Intelligence Review, 24:33{60}, 2005.

Stevens, Robert, Patricia Baker, Sean Bechhofer, Gary Ng, Alex Jacoby, Norman W. Paton, Carole A. Goble, and Andy Brass. TAMBIS: transparent access to multiple bioinformatics information sources. Bioinformatics 16(2):184-186, 2000.

Sean Bechhofer, Iain Buchan, David De Roure, Paolo Missier, John Ainsworth, Jiten Bhagat, Philip Couch, Don Cruickshank, Mark Delderfield, Ian Dunlop, Matthew Gamble, Danius Michaelides, Stuart Owen, David Newman, Shoaib Sufi, Carole Goble, Why linked data is not enough for scientists. Future Generation Computer Systems 29(2):599-611, 2013.

Compton, Michael, David Corsar, and Kerry Taylor. Sensor data provenance: SSNO and PROV-O together at last. Terra Cognita and Semantic Sensor Networks pp 67-82, 2014.