



D-SA3.2 – Scientific gateway: first report on release of services and management integration

18/10/2012

Project acronym: VERCE
Project n°: 283543
Funding Scheme: Combination of CP & CSA
Call Identifier: FP7-INFRASTRUCTURES-2011-2
WP: WP7/SA3, Scientific gateway: first report on release of services and management integration
Filename: D-SA3.2.pdf
Author(s): A. Spinuso (KNMI), L. Trani (KNMI), T. van Eck (KNMI)
Location: <http://www.verce.eu/Repository/Deliverables/RP2/>
Type of document: Deliverable
Dissemination level: Public
Status: Final
Due date of delivery: 1/10/12
Reviewer: Genevieve Moguilny, Siew Hoon Leong
Keywords: SA3, gateway, portal, components, registry.

<i>Version</i>	<i>Author</i>	<i>Date</i>	<i>Comments</i>
1	A. Spinuso KNMI	2/09/2012	Initial draft, basic layout.
2	A. Spinuso KNMI	7/09/2012	Scientific Gateway prototype section edited, to be reviewed
3	I. Klampanos UEDIN	19/09/2012	Description and cross reference to VERCE Registry
4	L. Trani KNMI	16/09/2012	Development Strategy and Cross project coordination edited, to be reviewed
5	A. Spinuso KNMI	21/09/2012	Executive Summary
5	A. Spinuso KNMI	24/09/2012	Reviewer's comment edited

Copyright notice

COPYRIGHT © VERCE PROJECT, 2011-2015. SEE www.verce.eu FOR DETAILS ON VERCE.

VERCE, *Virtual Earthquake and seismology Research Community e-science environment in Europe*, is a project co-funded by the European Commission as an Integrated Infrastructure Initiative within the 7th Framework Programme. VERCE began in October 2011 and will run for 4 years.

This work is licensed under the Creative Commons Attribution-Noncommercial 3.0 License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc/3.0> or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, and USA.

The work must be attributed by attaching the following reference to the copied elements:

COPYRIGHT © VERCE PROJECT, 2011-2015. SEE www.verce.eu FOR DETAILS ON VERCE. Using this document in a way and/or for purposes not foreseen in the license requires the prior written permission of the copyright holders. The information contained in this document represents the views of the copyright holders as of the date such views are published.

Contents

Executive Summary	4
1 Development strategy implementations	5
2 Cross project coordination (NERA, EPOS)	5
3 Scientific gateway prototype	5
3.1 Workflow implementation	6
3.2 Provenance visualisation	7
3.3 Deployment and first tests	7
3.4 VERCE Registry	10
4 Conclusions	10
Glossary	11
List of Figures	
1 A DISPEL function returning SeismoStage that performs the plotting of the data stream, storing the image in a specific location.	6
2 SeismoStage function and Provenance.	7
3 Gateway prototype: The Browser Based DISPEL language editor that allows the editing, validation, and submission of a workflow to the deployed platform.	8
4 Gateway prototype: provenance and data visualisation.	8
5 Gateway prototype: Deployment on DEP-UEDIN-01.	9

Executive Summary

The objective of the SA3 work package is to provide and operate a scientific gateway that enables the access to the services available through the VERCE platform. These services will allow the execution of a set of scientific applications that, accordingly to their requirements, will be implemented adopting the technologies indicated by JRA2.

In the last six months the work package effort has been focused, at the first stage, on the establishment of a development strategy for the implementation of a set of demonstrators of the scientific applications defined by the VERCE use cases. Moreover, attention has been also dedicated to the coordination with other on going projects in Seismology and Solid Earth Sciences, such as NERA (Network of European Research Infrastructures for Earthquake Risk Assessment) and EPOS (European Plate Observing System), in order to envisage how the VERCE technology could also provide valuable solutions to these other initiatives.

In this deliverable we present therefore the implementation of a working demonstrator of an initial scientific gateway. The main objective of this prototype is to show to the VERCE partners how the different layers of the platform interact and where are the boundaries between the development efforts of the IT experts and the scientists. The prototype demonstrates how a workflow defining a preprocessing pipeline of seismic traces can be implemented with the DISPEL language and how it translates into the execution of the scientist's analysis code within a distributed deployment of the platform on the DEP-UEDIN-01, the data Intensive machine in Edinburgh. The prototype shows how the data transformation, the collection of the metadata and its visualisation are performed in a complete automated way, leaving to the user only the important responsibility of validating the results. A set of interactive web pages for the workflow submission and provenance visualisation is also provided, as part of the initial design of the Scientific gateway front end.

The demonstrator definitely fostered the discussion on how to proceed with the work in the following months. In that respect, relevant contacts have been established among the scientists and the IT experts, in order to form small working groups that will have the responsibility to push forward the development of the scientific use cases. The VERCE scientific gateway's components will strictly follow these developments in order to refine its specification, towards an incremental implementation of the users' needs and expectations.

1 Development strategy implementations

Given the ambitious goals of the project and consequently the risks coupled to them, we decided to adopt an implementation plan constituted by short cycles of about 6 months. During these cycles we aim to deliver incrementally working components of the infrastructure, starting from the specification of two main use cases: a data intensive use case and a compute intensive use case. The main actors involved in the use cases have been identified to establish a cross WP collaboration in order to foster a strong targeted interaction among the partners and facilitate the developments. The initial cycle envisages specific functionalities to be implemented. On the data intensive side the developments go in the direction of a better integration of the scientific code in the infrastructure and in general the set up of a full preprocessing pipeline, including the initial definition of the user interaction. Instead regarding the compute intensive use case we aim at integrating mechanisms to interface HPC infrastructures within the current framework, including data staging functionalities, starting with the installation and testing of technologies to ship data around the identified testbeds. In the coming months we will continue with the technological investigation towards solutions that could better suit to the requirements, and this activity will be carried out with the strict collaboration of JRA2.

2 Cross project coordination (NERA, EPOS)

The coordination with other seismological projects is strong and guaranteed by the engagement and full integration with the IT developments currently ongoing in the other similar initiatives. As an example we can mention NERA¹ (Network of European Research Infrastructures for Earthquake Risk Assessment) which is currently looking at the possibility to adopt the VERCE technology to produce and provide to users products adopting ad hoc workflows. Another important aspect is the possible integration of derived products like "synthetics" computed on the VERCE platform within existing and established datacenters for long term preservation and curation. A particular attention is also given to the initiatives addressing the definitions of data formats and "standards" in the seismological field, like the ones currently carried on by EPOS² (European Plate Observing System) and NERA. VERCE will keep a close collaboration to insure the best compatibility and interoperability as possible. As for the data access, the storage and curation of intermediate and final data products, it is important to adopt technologies and formats already used in larger context avoiding temporary ad hoc solutions. If from one side this can result in a more complex and long process with all the related risks mainly due to different cross projects alignments, in the long run it seems to be the only way to go to ensure maintainability. In order to support this strategy important collaborations have been established with key people and groups across the projects.

3 Scientific gateway prototype

The development of a first prototype of the scientific gateway has been conducted accordingly to the objectives defined within the Data Intensive Use Case. The use case requires the implementation of a preprocessing pipeline for seismic traces, as the first phase of the cross correlation analysis. An initial release and deployment has been provided and demonstrated during the first VERCE training meeting, held in Liverpool on 3-4 September 2012. In the following sections we describe the core components of the prototype, its deployment environment and preliminary tests. Moreover a registry of processing elements has been also developed and shown to the user. We'll provide a brief introduction of its user interface and functionalities.

¹<http://www.nera-eu.org/>

²<http://www.epos-eu.org/>

```

1 <SeismoStage> plot (String plotLocation,
2                     String provenanceRes,
3                     String preprocessRes)
4 {
5
6
7     SeismoMetadataTuple metastore = new SeismoMetadataTuple;
8     WaveformPlot plot=new WaveformPlot;
9
10    |-provenanceRes-|=>metastore.resource;
11    |- repeat enough of REQUEST_ID -|=>metastore.processid;
12    |-preprocessRes-|=>plot.resource;
13    |-repeat enough of ["filedestination="+plotLocation]-|=>plot.parameters;
14
15    plot.metadata=>metastore.metastring;
16
17    return SeismoStage (<Connection input=plot.input;
18                        Connection stepbackid=metastoreexl.stepbackid>=>
19                        <Connection datasetid=metastore.datasetid;
20                        Connection output=plot.output;
21                        Connection metadata=metastore.processedMetadata>);
22
23 };

```

Figure 1: A DISPEL function returning SeismoStage that performs the plotting of the data stream, storing the image in a specific location.

3.1 Workflow implementation

The work has been conducted starting from an early implementation which was developed and demonstrated for the ADMIRE project. The workflow performs a sequence of data transformations on a number of seismic traces extracted from a waveform data archive. In the last six months it has been further improved, in collaboration with JRA2, accordingly to the use case requirements.

The workflow consists in a DISPEL script that uses seismological processing elements (PE) within high level DISPEL functions (Figure 1). Each of these functions accepts several configuration parameters and returns a new composite PE, which can be inserted within a pipeline and eventually executed within a distributed environment hosting the platform.

Besides the transformations on the data stream, these composite PEs allow the production and the storage of provenance information in a centralised provenance database, independent from the distribution of the execution (Figure 2).

This new implementation offers a better integration of OGSA-DAI and Python/Obspy³, which has been updated to support the streaming of SEED⁴ data. The DISPEL code has been reengineered in order to improve its scalability and the storage of metadata, exploiting as much as possible the prerogative of DISPEL as a workflow composition language for data-intensive applications.

To run the workflow we are adopting a browser based interface that allows the user to type in the DISPEL code, to validate it and submit it to the deployed infrastructure (Figure 3). This is part of the initial front end of the VERCE scientific gateway, that will be later on substituted by a more interactive interface where the DISPEL code will be hidden to the end user's eyes. Accordingly with the customisable sections of the workflow, we foresee that the user will be exposed to a webpage hosting interactive

³<http://obspy.org/>

⁴http://www.iris.edu/manuals/SEED_chpt1.htm

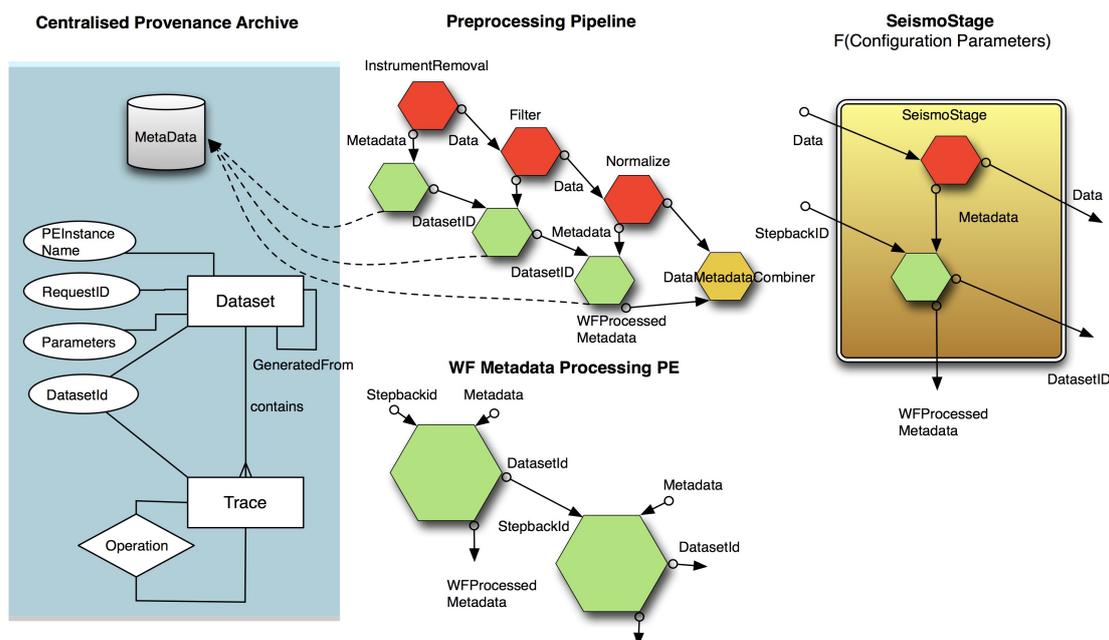


Figure 2: SeismoStage function and Provenance.

widgets such as forms, interactive maps and simple visual pipelines. The latter can be used, for example, to compose in different ways a certain number of DISPEL SeismoStages.

3.2 Provenance visualisation

As already explained in the previous section, the prototype is capable of storing the metadata produced by the transformations applied on the data stream. This information has to be properly visualised in order to give immediate feedback to the user on the current status of the computation, in order to support the validation of the results and the diagnostic of errors, if occurred.

As initial view into this wealth of information, we proposed to the users a browser based web interface where they can browse through the data navigating among three different pages (Figure 4).

Process view: Near real time monitoring of all the processing metadata relative to each transformation applied to each piece of the data stream. Some of the most important metadata displayed here are the following: a series of IDs identifying the processing element, the dataset produced and the machine where the computation was carried out. The page also shows the parameters used within the computation and the creation timestamp of the metadata. Moreover, whether the workflow execution produced downloadable datasets and visualisation files, a link pointing to these resources is also provided.

Provenance Trace view: Same set of metadata as the Process view. This page provides a view of the full history of the transformations applied to a single segment of the data stream.

Dataset Detail view: This page provides to the user the detailed information about the status of the content of the dataset, at a certain point of the pipeline.

3.3 Deployment and first tests

The prototype has been deployed on four nodes of the DEP-UEDIN-01 machine hosted by the University of Edinburgh. The set up of the prototype and the distribution of the workflow execution is shown in

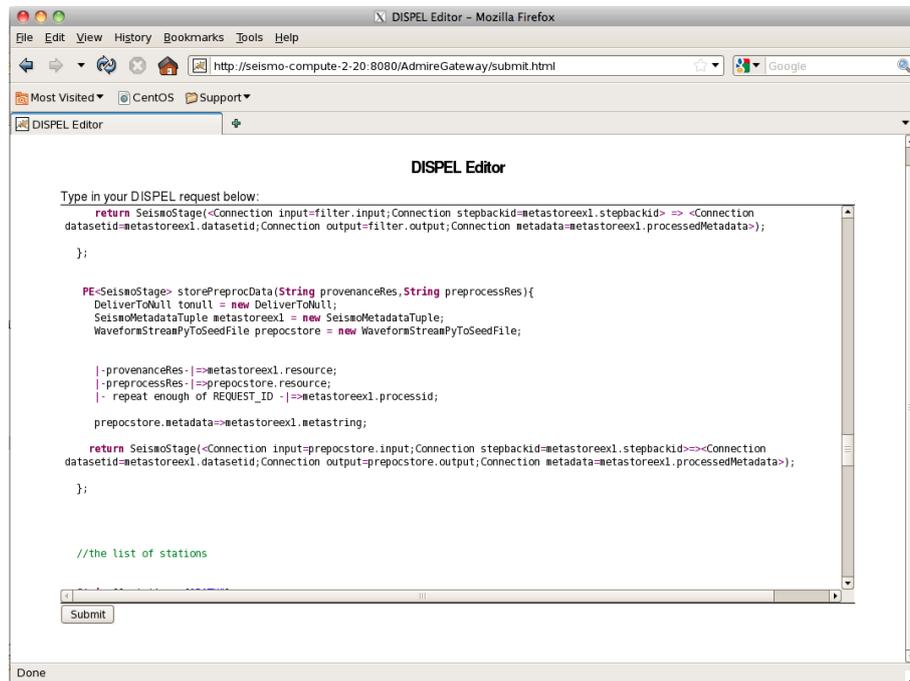


Figure 3: Gateway prototype: The Browser Based DISPEL language editor that allows the editing, validation, and submission of a workflow to the deployed platform.

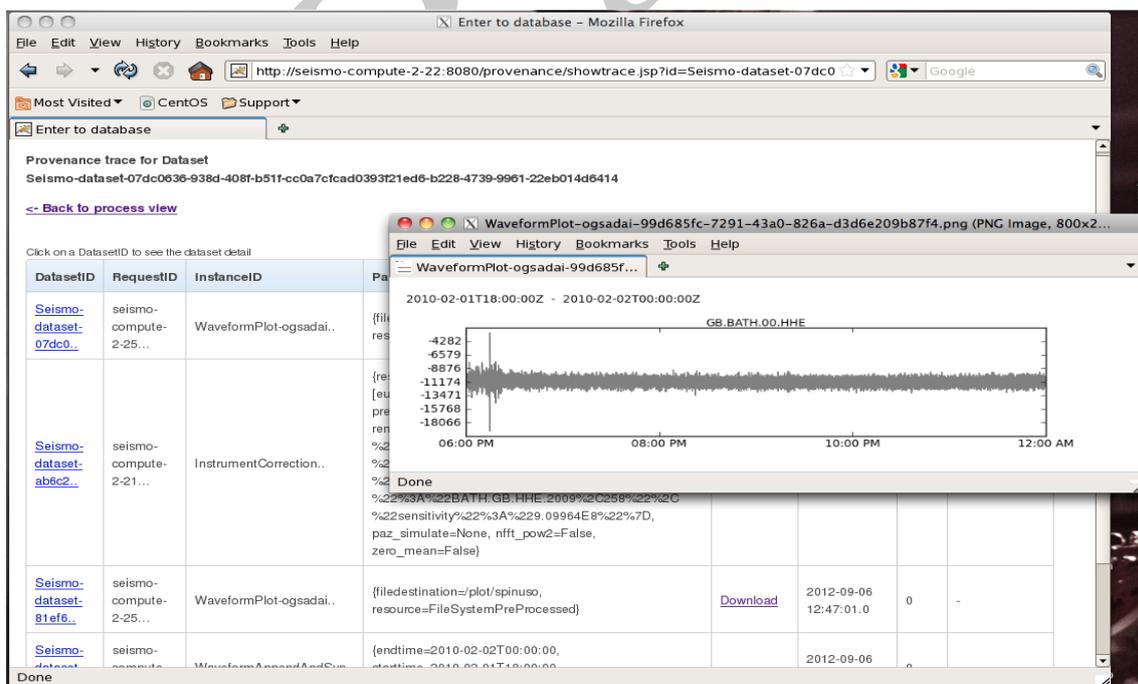


Figure 4: Gateway prototype: provenance and data visualisation.

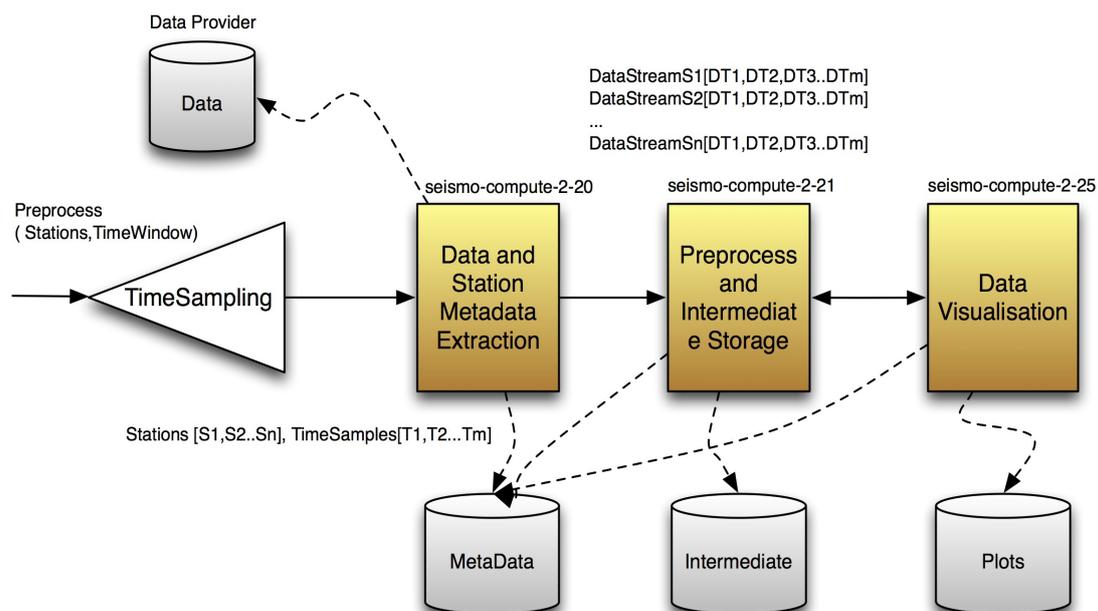


Figure 5: Gateway prototype: Deployment on DEP-UEDIN-01.

Figure 5.

Each node uses a 2-core Intel Atom CPU (1.6 GHz), an nVidia Ion GPU and 4GB of DDR3 memory. It contains a single 256GB solid state drive upon which to store appliance and temporary data, as well as three 2TB hard disk drives for permanent data storage.

In our test we stream 15 days of data belonging to a single station, extracted and synchronised from a local archive (~140 mb). The data is converted into a stream composed by segments of 12 hours of continuous waveform time-series. The preprocessing pipeline consists in the following sequence of six stages:

1. **Visualisation:** The plot of the raw data is generated and made available to the user.
2. **Instrument Correction:** A correction of the measurements eroding the effect of the instrument on the recorded signal.
3. **Filtering:** A band-pass filter that passes frequencies within a certain range and rejects (attenuates) frequencies outside that range.
4. **Whitening:** A process to improve the resolution and appearance of seismic data through a crude attempt to correct for frequency attenuation.
5. **Storage:** Each segment of the stream is stored and made available for downloading.
6. **Visualisation:** The final plot gets generated and made available to the user.

Provenance metadata is recorded for each of these steps in a centralised database. These information are made available to the users via a tabular browser based interface, as already described in the previous section.

The execution of the entire workflow with the current setup completes in ~20 minutes. This results are now subjected to evaluation by the seismologists that will contribute in validating, correcting and extending the library of analysis algorithms already in use. Moreover a set of controlled benchmarks

will be performed to obtain more quantitative results that will contribute to make a proper evaluation on the potential of the VERCE platform.

3.4 VERCE Registry

The VERCE registry is designed to allow the registration of processing elements and other workflow components, such as *functions*, *types*, etc. It provides a RESTful interface to the Workflow engine components of VERCE and in particular to the Dispel Gateways. This interface allows these components to have a common, uniform and up-to-date view of the workflow ecosystem available throughout the VERCE platform. The VERCE Registry is described in more detailed in Section 2.1.3 of D-JRA2.1.1.

4 Conclusions

In the last months a major effort has been dedicated to the implementation of the seismic preprocessing pipeline, adopting the technologies provided by JRA2. This activity was extremely important to start a fruitful discussion on how to proceed with the development of the use cases, which led to the composition of small and heterogeneous working groups. In that respect a series of meetings is already planned for the upcoming months, where the seismologist and the IT expert will work together to tackle the implementation challenges of the use cases. They will identify with higher granularity all those relevant details that will also impact on the interactive features of the VERCE scientific gateway.

DRAFT

Glossary

component One of the computational elements involved in a data-intensive or computational process, such as: application codes, scripts, workflows, services, catalogues, registries, data collections, data resources, functions, gateways, libraries, PEs, PE instances, format definitions and types.

data archive The long-term storage of scientific data and methods.

data-intensive An adjectival phrase that denotes that the item to which it is applied requires attention to the properties of data and to the ways in which data are handled.

Dispel Data-Intensive Systems Process Engineering Language, a workflow composition language for data-intensive applications.

gateway A software subsystem, typically at the middleware level, that accepts requests for computational and data-handling tasks. It vets those requests to establish whether they are valid, e.g. are syntactically and semantically consistent, and are authorised. Requests that are not validated are rejected. Requests that are accepted are passed to other software systems, at the same or other locations, for execution. The gateway may partition and translate requests in order to combine heterogeneous services.

high-performance computing (HPC) Use of powerful processors, high-speed networks and parallel supercomputers for running computationally intensive applications.

metadata Data that describes data. Metadata may include references to schemas, provenance, and information quality. In Seismology, metadata may also refer to data required in order to sanitise a seismograph's response.

OGSA-DAI Open Grid Service Architecture Data Access and Integration, an open source product for distributed data access and management.

ORFEUS Observatories and Research Facilities for European Seismology.

portal In the context of knowledge discovery, a tool designed for a particular group of domain experts that can be used via their browsers; it enables them to establish their identity and rights, and to pursue conveniently a set of research tasks for which the portal is designed.

processing element – PE A software component that encapsulates a particular functionality and can be used to construct a workflow.

registry A persistent store of definitions and descriptions of data or software components and their relationships accessed by tools and other elements of a distributed research environment. It is intended to facilitate discovery and use of the components.

repository A store holding software definitions, other shared code and data, that supports distributed concurrent access, update and version management.

science gateway A consistently presented set of facilities designed to be a convenient working environment for researchers in a particular domain, in this case seismology. It should bring together access to all of the capabilities and resources such a researcher needs: including catalogues of available data and tools, established methods and arrangements for applying them with specified parameters to specified data.

VERCE e-Infrastructure An envisaged result of VERCE, as an integrated computational and data environment that presents a coherent virtual research environment in which to conduct seismology research and eventually research in other Earth sciences.

VERCE Platform The current realisation of the VERCE e-Infrastructure at any time in the VERCE project. Initially this is not fully integrated and may only constitute a partial implementation. Nev-

ertheless, it is sufficient both to pursue research identified as priority seismology use cases and to develop and test the design of the VERCE e-Infrastructure. The VERCE Platform is an approximation to the VERCE e-Infrastructure. These approximations should converge on the VERCE e-Infrastructure by the end of the VERCE project. virtual research environment (VRE) A presentation of (ideally all of) the resources a researcher may need in a consistent and easily used form. These resources include catalogues, data, metadata, libraries, tools, workflows, programs, services, visualisation systems and research methods.

workflow A process of composed data-handling tasks, computational tasks and human interactions intended to implement a research method or established working practice.

DRAFT