



Joining B2SAFE

About

Document for Community Data Managers and Site Administrators that describes how communities can become part of the B2SAFE replication network in EUDAT.

Modified: 29 January 2018

Synopsis

This document is targeted at communities which want to deploy B2SAFE and become a node in the B2SAFE replication network. It explains the functionality that B2SAFE offers in the “joining” mode, the deployment workflow and which software needs to be installed.

Acronyms

PID: Persistent identifier associated to a file or iRODS collection, usually an EUDAT Handle.

ROR: Repository of Records, (persistent) identifier to the original data object. Can be of any identifier type. If the community does not assign an own identifier the B2SAFE PID will be used.

FIO: First ingested object, the persistent identifier associated to the very first object in in the EUDAT domain. If the the chain has only two elements, the master copy and the first replica, then the PARENT = FIO.

PARENT: Parent PID, the persistent identifier associated to the source object in a replication chain. If the chain has only two elements, the master copy and the first replica, then the PARENT = FIO.

Digital entity: Files and folders

Digital object: Files and folders that carry a persistent identifier and possibly some metadata.

Introduction

B2SAFE offers safe data replication across different data centres. Communities, repositories and data projects can use B2SAFE to distribute valuable data across the EUDAT network in order to keep it safe and to bring it closer to compute infrastructure. In the rest of this section we explain how B2SAFE works.

The underlying assumption of the safe replication of B2SAFE is that the data which needs to be replicated is stored in a repository, the so-called Repository of Records (ROR). In general it is assumed that the data in the repository of records is not to change, i.e. new files may be added to the repository but files do not change over time and are not deleted. The single data files in the repository of records receive a persistent identifier (PID) which is used to link the original data and the replicas. Moreover, the PID is used to store information for integrity checks, e.g. checksums. The ROR is used as a reference to test the integrity between original files and their replicas. This implies that in general we assume that the ROR does not change.

Upon replication from the ROR, new files at other EUDAT sites are created and need to be linked to their respective parent or to their respective originals. To this end B2SAFE makes use of specific PID entries: *ROR*, *FIO* and the parent PID (*PARENT*). The value in the field *ROR* designates the identifier of the data file in the

EUDAT receives funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement No. 654065.

community repository of records, while the value in the field *FIO* is the PID of the very first ingested in the EUDAT domain. The object's *PARENT* is the PID of the direct parent of the replica. This ensures that we can retrieve the original files of each replica and have reference data objects for integrity checks. To find replicas when given the PID of the original file, the PID field *REPLICA* is used, where all PIDs of direct replicas of the file are stored.

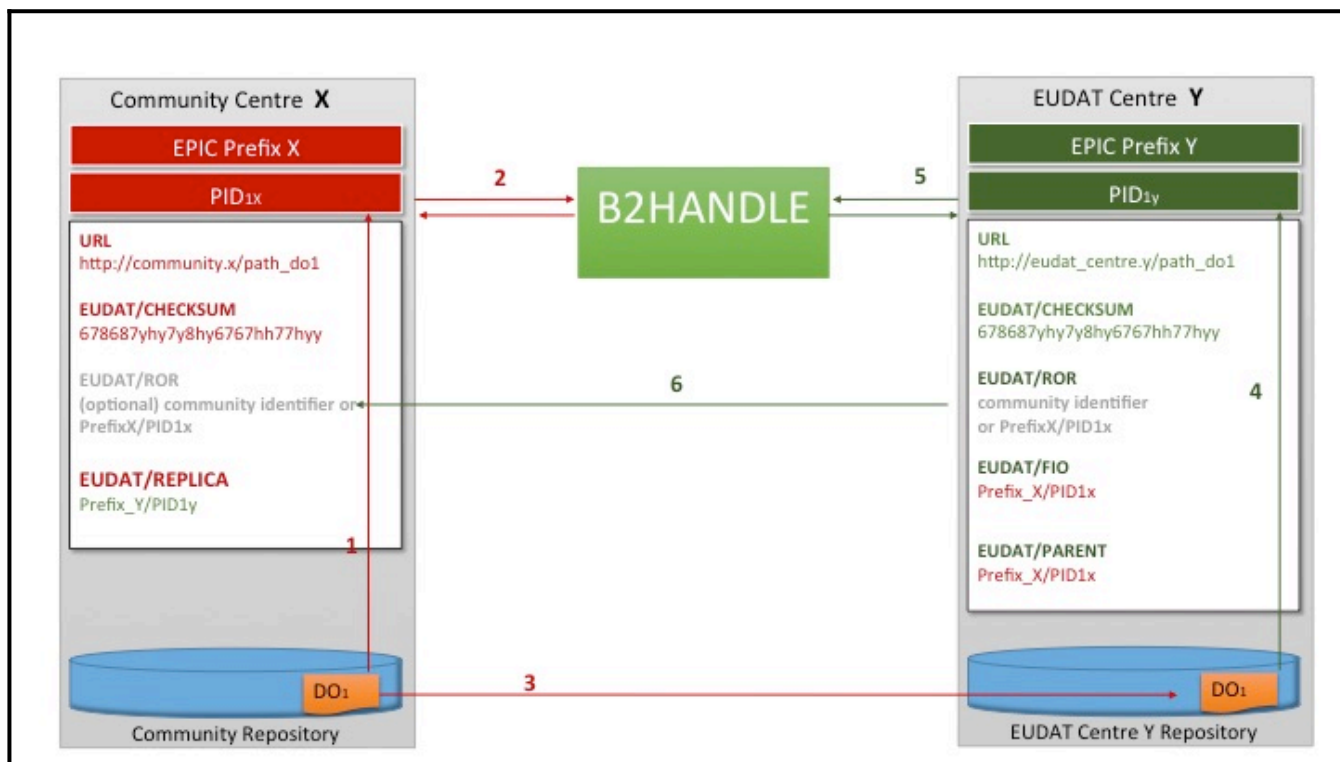


Figure 1: Linking replicas to their original data with PIDs

Figure 1 depicts the flow of linking replicas with their original data. The flow is as follows:

1. The Data Object (DO) is stored at a centre. The figure shows a Community Centre and the workflow discusses it as such, but this is not a requirement.
2. The (community) centre labels the data with a PID which holds information on the data location (URL), its original location (as URL) and some additional data such as the checksum.
3. The data is replicated with B2SAFE to an EUDAT centre Y.
4. Upon replication, the copy of the data object receives another PID. This new PID contains an additional entry ROR, which links to the DO at the (community) centre, a link to the FIO (here we assume the community centre is part of the EUDAT network) and a link to the direct parent, which in this case is the same as the FIO.
5. The PID for the replica is created under a different prefix than the PID for the original DO.
6. The PID of the replica is entered under *REPLICA* in the PID of the original DO.

The replication of files, the generation of PIDs and the linking between originals and replicas are automated by iRODS workflows written as iRODS rules; see [the B2SAFE user documentation](#) for some example workflows. For more information on the PID linking in B2SAFE also please refer to [the B2SAFE documentation](#).

As a centre offering B2SAFE to communities you will become part of the EUDAT data replication network, allowing data to be replicated from you to other centres or accepting data from other data and community centres via the



replication chain. To allow for the latter, data centres need to declare how much storage they would like to offer.

General workflow

The following workflow applies in order to join B2SAFE.

1. Deployment of an iRODS/B2SAFE instance (for technical requirements see below)
2. Agreements on iRODS federations with other EUDAT centres and community centres
3. Entry in the [Resource Coordination Tool](#) (RCT) registry making the new B2SAFE node known to EUDAT
4. (Optional) Entry in the RCT registry offering a certain amount of storage
5. (Optional) Deployment of B2STAGE, offering a gridFTP endpoint to access data in B2SAFE

Technical requirements

A community wanting to join B2SAFE will need to deploy the following software:

1. iRODS 4.1.10 or 4.2.1 server
2. Deployment of the B2SAFE software stack
3. Installing the [B2HANDLE](#) Python library and having access to a Handle prefix to create and manage PIDs

Support

Support for BSAFE is available via the EUDAT ticketing system through the [webform](#).

If you have comments on this page, please also submit them through the [EUDAT ticketing system](#).

Document Data

Version: 2.1

Authors:

Christine Staiger, christine.staiger@surfsara.nl

Giovanni Morelli, g.morelli@cineca.it

Editors:

Themis Zamani, themis@grnet.gr

Kostas Kavoussanakis k.kavoussanakis@epcc.ed.ac.uk

[Read more](#)