



#### **D.4.4.2 SIP-AIP conversion component: Reference Implementation**

**DOI: 10.5281/zenodo.1172925**

|                         |  |
|-------------------------|--|
| Grant Agreement Number: | 620998   |
| Project Title:          | European Archival Records and Knowledge Preservation |
| Release Date:           | 14 <sup>th</sup> February 2018                       |
| <b>Contributors</b>     |  |
| <b>Name</b>             | <b>Affiliation</b>                                   |
| Sven Schlarb            | Austrian Institute of Technology                     |
| Jan Rörden              | Austrian Institute of Technology                     |
| Mihai Bartha            | Austrian Institute of Technology                     |
| Kuldar Aas              | National Archives of Estonia                         |
| Andrew Wilson           | University of Brighton                               |
| David Anderson          | University of Brighton                               |
| Janet Anderson          | University of Brighton                               |

## **STATEMENT OF ORIGINALITY**

**Statement of originality:**

This deliverable contains original unpublished work except where clearly indicated otherwise. Acknowledgement of previously published material and of the work of others has been made through appropriate citation, quotation or both.

## TABLE OF CONTENTS

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Scope of this document</b> .....          | <b>6</b>  |
| <b>2</b> | <b>Relation to other documents</b> .....     | <b>6</b>  |
| <b>3</b> | <b>Introduction</b> .....                    | <b>6</b>  |
| <b>4</b> | <b>E-ARK Web Architecture overview</b> ..... | <b>8</b>  |
| 4.1      | Cluster software stack .....                 | 8         |
| 4.2      | Standalone software stack .....              | 10        |
| <b>5</b> | <b>SIP-AIP conversion component</b> .....    | <b>15</b> |
| 5.1      | Task execution framework .....               | 15        |
| 5.2      | SIP-AIP conversion tasks .....               | 16        |
| <b>6</b> | <b>E-ARK Web user guide</b> .....            | <b>19</b> |
| 6.1      | SIP Creator .....                            | 19        |
| 6.1.1    | Initialise SIP package creation.....         | 19        |
| 6.1.2    | Adding files to the SIP .....                | 19        |
| 6.1.3    | SIP creation process .....                   | 22        |
| 6.2      | SIP to AIP conversion .....                  | 25        |
| 6.2.1    | Start SIP to AIP conversion process .....    | 25        |

## List of Figures

|  |    |
|--|----|
| Figure 1: The scope of the this document is Deliverable D4.4 Part B .....                                  | 5  |
| Figure 2: Screenshot of the initial page of the web application E-ARK Web .....                            | 7  |
| Figure 3: Architecture overview of the full-scale version of E-ARK Web .....                               | 9  |
| Figure 4: Architecture overview of the lightweight version of E-ARK Web.....                               | 10 |
| Figure 1 Docker containers of the Standalone Software Stack and the main links between the containers..... | 13 |
| Figure 2 Origin of the images which are the basis to run the Docker containers .....                       | 14 |
| Figure 3 Data directories of the Docker containers and mapping to directories of the host file system..... | 15 |
| Figure 8: Parent-child aggregation submission process.....   | 32 |

## List of Tables

|  |    |
|--|----|
| Table 1 Docker containers of the standalone deployment stack ..... | 11 |
| Table 2 SIP-AIP conversion tasks .....                             | 18 |

## Executive Summary

This deliverable D4.4 Part B is the second part of deliverable D4.4 and describes the SIP-AIP conversion component in its final version. The first part D4.4 Part A is an update of the previous deliverable D4.3 and represents the final version of the AIP specification. Figure 1 illustrates the two parts A and B of deliverable 4.4 and highlights Part B as the scope of the present document.

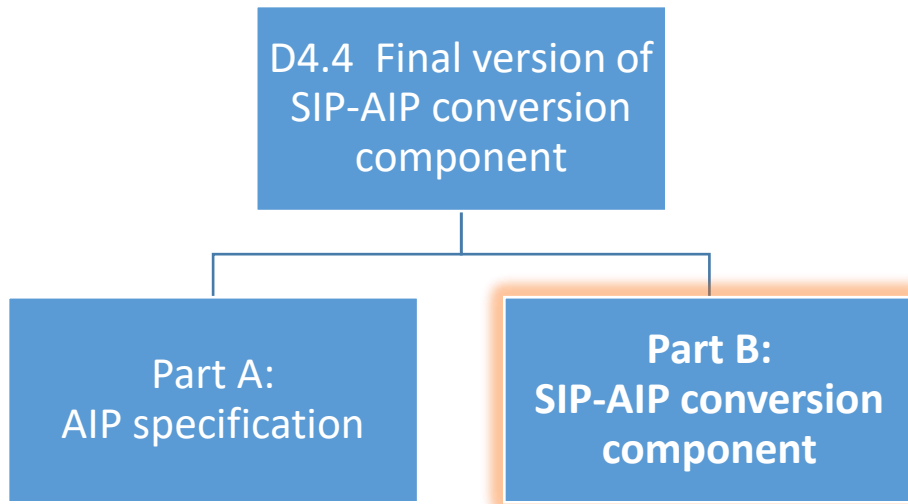


Figure 1: The scope of the this document is Deliverable D4.4 Part B

This document focuses on the E-ARK reference implementation of the SIP-AIP conversion component and describes the implementation of the basic concepts and definitions. The deliverable is a follow-up version of E-ARK deliverable D4.3 “E-ARK AIP pilot specification”.<sup>1</sup>

The SIP-AIP conversion component is the reference implementation of the AIP format specification and is an integral part of the web application E-ARK Web (or in short: *earkweb*).<sup>2</sup> The SIP-AIP conversion component consists of a set of individual tasks which are executed in a specific order to convert an E-ARK submission information package (SIP) into the E-ARK archival information package (AIP).

---

<sup>1</sup> <http://eark-project.com/resources/project-deliverables/53-d43earkaipspec-1>

<sup>2</sup> <https://github.com/eark-project/earkweb>; the actual SIP to AIP conversion component consists of a set of tasks in <https://github.com/eark-project/earkweb/blob/master/workers/tasks.py>.

## 1 Scope of this document

This document describes the SIP-AIP conversion component as part of the reference implementation E-ARK Web (in short: *earkweb*.)<sup>3</sup> It describes a set of tasks which is used to convert an E-ARK SIP to an E-ARK AIP and explains how these tasks are employed by the web application frontend of *earkweb*. In order to explain the context in which the SIP-AIP component is embedded, some details about the architecture and a user guide to the *earkweb* reference implementation are given.

## 2 Relation to other documents

This document describes how the SIP-AIP conversion is implemented as part of the component *earkweb*.<sup>4</sup> It complements part A “AIP format specification” of deliverable D4.4, “Final version of SIP-AIP conversion component”.

Deliverable D4.4 has two previous iterations, Deliverable D4.2<sup>5</sup> and Deliverable D4.3.<sup>6</sup>

The SIP-AIP conversion component was designed to integrate with the scalable, Hadoop-based back-end implementation. The deliverable D6.3 “Data Mining Showcase” includes a detailed overview of the updated architecture of the Integrated Platform Reference Implementation (IPRIIP). The SIP specification is described in the D3.x deliverables.

The fundamental document to understand the purpose of the SIP-AIP conversion is the Reference Model for an Open Archival Information System (OAIS).<sup>7</sup>

## 3 Introduction

*earkweb* is an open source archiving and digital preservation system. It is OAIS<sup>8</sup>-oriented which means that data ingest, archiving and dissemination functions operate on information packages bundling content and

---

<sup>3</sup> <https://github.com/eark-project/earkweb>; the actual SIP to AIP conversion component consists of a set of tasks in <https://github.com/eark-project/earkweb/blob/master/workers/tasks.py>.

<sup>4</sup> <https://github.com/eark-project/earkweb>; the actual SIP to AIP conversion component consists of a set of tasks in <https://github.com/eark-project/earkweb/blob/master/workers/tasks.py>.

<sup>5</sup> <http://eark-project.com/resources/project-deliverables/18-d42-e-ark-aip-draft-specification>

<sup>6</sup> <http://eark-project.com/resources/project-deliverables/53-d43earkaipspec-1>

<sup>7</sup> Reference Model for an Open Archival Information System (OAIS); Retrieved from <http://public.ccsds.org/publications/archive/650x0m2.pdf>, in the following we are always referring to this version 2 of the OAIS which was published in June 2012 by CCSDS as “magenta book” (ISO 14721:2012). The document is cited using the abbreviation OAIS.

<sup>8</sup> <http://public.ccsds.org/publications/archive/650x0m2.pdf>

metadata in contiguous containers. The information package format uses METS<sup>9</sup> to represent the structure and PREMIS<sup>10</sup> to record digital provenance information.

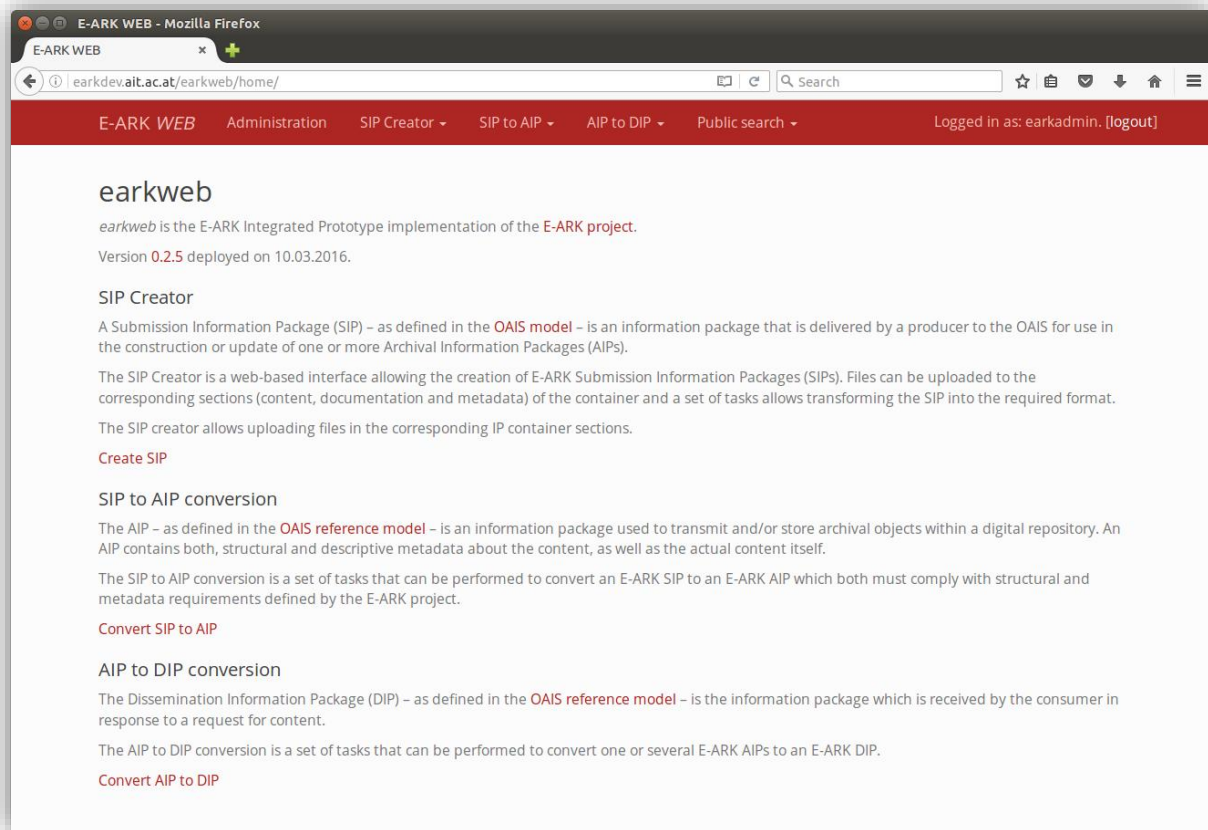


Figure 2: Screenshot of the initial page of the web application E-ARK Web

*Earkweb* offers functionality for the three types of information packages defined in the OAIS reference model: the Submission Information Package (SIP) which is the information sent from the producer to the archive, the Archival Information Package (AIP) which is the information stored by the archive, and the Dissemination Information Package (DIP) which is the information sent to a user when requested. The system allows executing different types of actions, such as information extraction, validation, or transformation operations, on information packages to support ingesting a SIP, archiving an AIP, and creating a DIP from a set of AIPs.

---

<sup>9</sup> <http://www.loc.gov/standards/mets>

<sup>10</sup> <http://www.loc.gov/standards/premis>

*Earkweb* consists of a frontend web application together with a task execution system based on Celery<sup>11</sup> which allows synchronous and asynchronous processing of information packages by means of processing units which are called “tasks”.

The backend can also be controlled via remote command execution without using the web frontend. The outcomes of operations performed by a task are stored immediately so that the status information in the frontend's database can be updated afterwards.

The SIP to AIP conversion component is implemented as a set of backend Celery tasks. These tasks can be executed using the web application frontend, by invoking the tasks in headless mode. *Earkweb* also offers a pre-defined workflow for batch processing which executes the full chain of tasks for automatic SIP-AIP conversion. These alternatives will be explained in detail throughout this document.

## 4 E-ARK Web Architecture overview

### 4.1 Cluster software stack

The E-ARK Web architecture is designed for efficiently processing, storing, and accessing very large data collections in terms of scalability, reliability, and cost. The system makes use of technologies like the Apache Hadoop<sup>12</sup> framework, NGDATA's Lily<sup>13</sup> repository, and the Apache Solr<sup>14</sup> search server allowing the repository infrastructure to scale-out horizontally. Using Hadoop, the number of nodes in a cluster is virtually unlimited and clusters may range from single node installations to clusters comprising thousands of computers. The diagram in Figure 3 gives an overview about this architecture.

---

<sup>11</sup> <http://www.celeryproject.org/>

<sup>12</sup> <http://hadoop.apache.org/>

<sup>13</sup> <https://github.com/NGDATA/lilyproject>

<sup>14</sup> <http://lucene.apache.org/solr/>



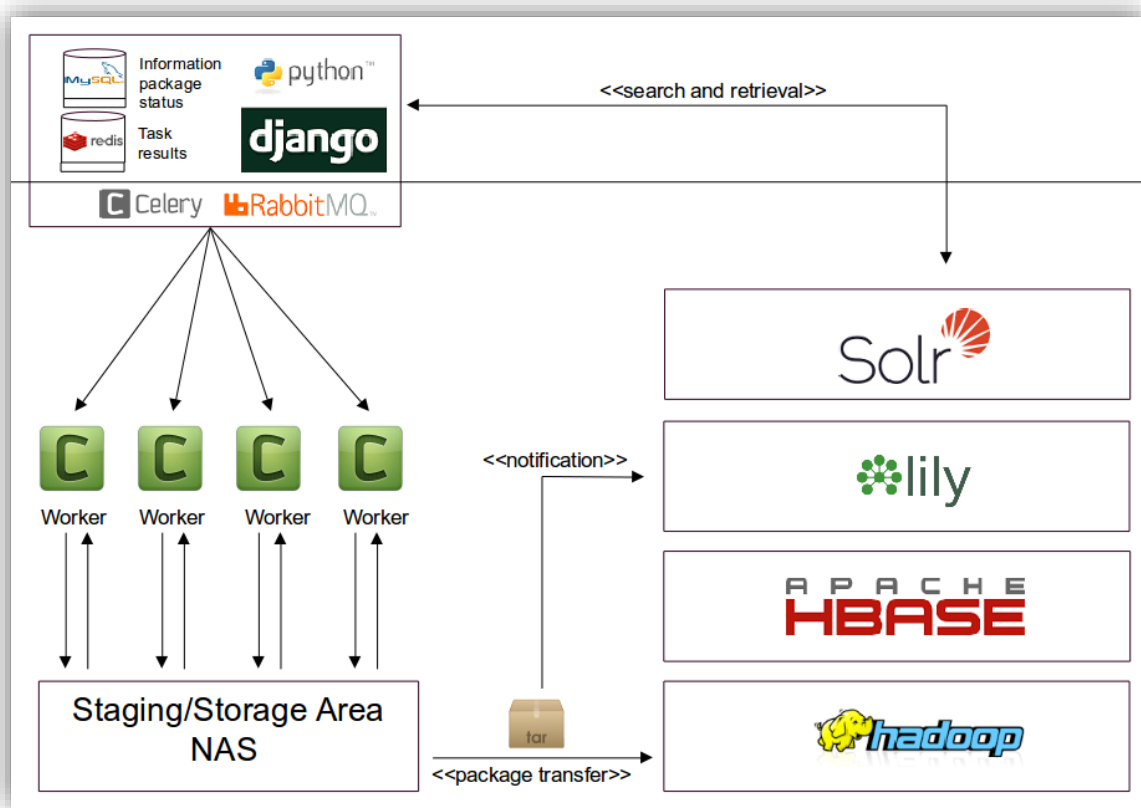


Figure 3: Architecture overview of the full-scale version of E-ARK Web

The user interface, represented by the box on top of the diagram, is a Python/Django-based web application which allows managing the creation and transformation of information packages. It supports the complete archival package transformation pipeline, beginning with the creation of the Submission Information Package (SIP), to the conversion to an Archival Information Package (AIP), to the creation of the Dissemination Information Package (DIP) which is used to disseminate digital objects to the requesting user. Tasks can be assigned to Celery workers (green boxes with a "C") which share the same storage area and the result of the package transformation is stored in the working directory.

Once the creation of information packages is finished, they can be deployed to the Lily access repository. Lily is built on top of HBase,<sup>15</sup> a NoSQL database that is running on top of Hadoop. Lily defines some data types most of which are based on existing Java data types. Lily records are defined using these data types rather

<sup>15</sup> <https://hbase.apache.org/>

than using plain HBase tables, which makes them better suited for indexing due to a richer data model. The Lily Indexer is the component which sends the data to the Solr server and keeps the index synchronized with the Lily repository. Solr neither reads data from HDFS<sup>16</sup> nor writes data to HDFS. The index is stored on the local file system and optionally distributed over multiple cluster nodes if index sharding (mechanism used to spread load on databases across multiple servers) or replication is used.

## 4.2 Standalone software stack

There is also a lightweight version of E-ARK Web where the large-scale storage backend (HDFS, HBase) is replaced by a conventional file system storage and the Solr search server is a single instance of Solr instead of a Solr Cloud deployment, as illustrated by the diagram in Figure 4.

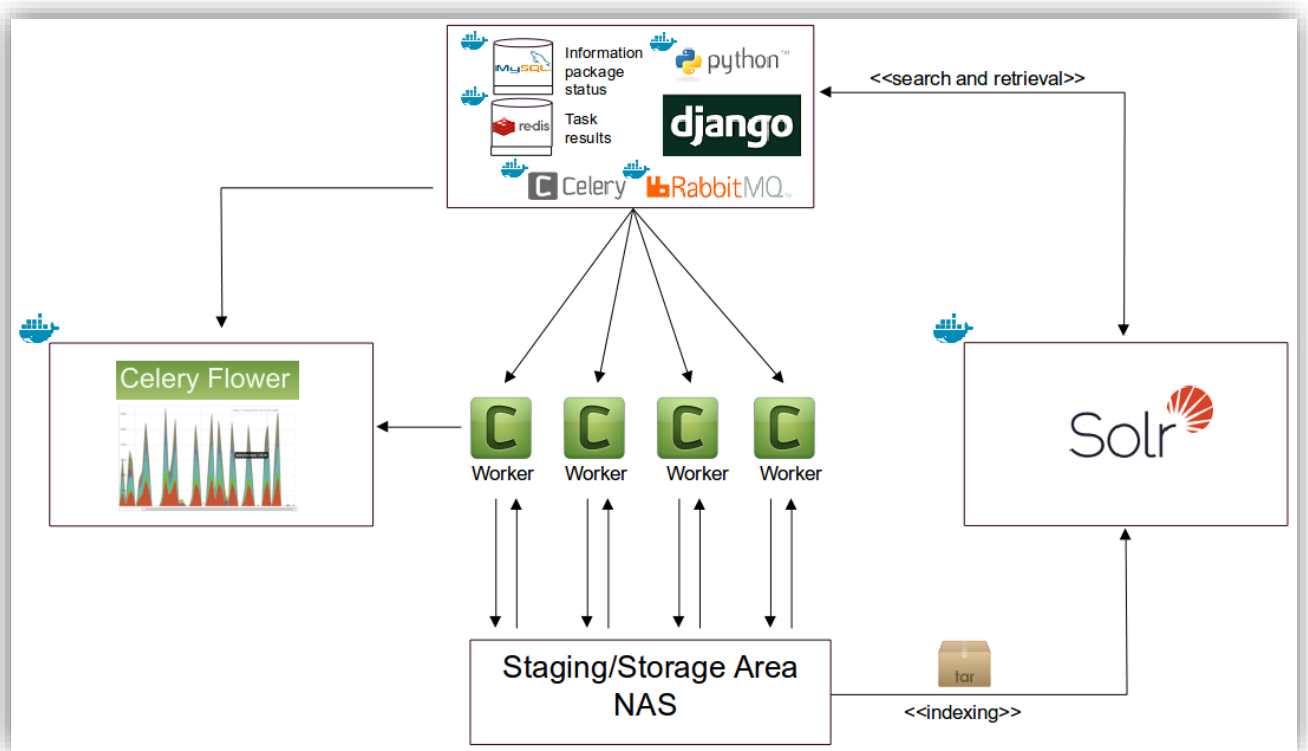


Figure 4: Architecture overview of the lightweight version of E-ARK Web

The standalone version of the IPRIP e-archiving environment has been made available using a container-based

<sup>16</sup> „Hadoop Distributed File System (HDFS) is designed to reliably store very large files across machines in a large cluster.”, <https://wiki.apache.org/hadoop/HDFS>

deployment model based on Docker<sup>17</sup> in order to support simple and modular installation of the software.

Docker is an open-source engine that automates the deployment of any application as a lightweight and portable container that will run on any platform where the Docker engine is supported.<sup>18</sup> In order to allow deploying IPRIP services on a Docker platform, Docker containers for the individual services of *earkweb*'s frontend and backend have been created.

Table 1 shows the software module and the corresponding image which is used in a Docker deployment. The left column designates the software module and the right column indicates the image which is used to create and run a container to provide the corresponding service.

| Software module             | Docker image                 |
|-----------------------------|------------------------------|
| MySQL <sup>19</sup>         | earkdbimg <sup>20</sup>      |
| Solr <sup>21</sup>          | Solr <sup>22</sup>           |
| RabbitMQ <sup>23</sup>      | tutum/rabbitmq <sup>24</sup> |
| Redis <sup>25</sup>         | tutum/redis <sup>26</sup>    |
| E-ARK Web <sup>27</sup>     | earkwebimg <sup>28</sup>     |
| Celery <sup>29</sup>        | earkwebimg                   |
| Celery Flower <sup>30</sup> | earkwebimg                   |

Table 1 Docker containers of the standalone deployment stack

Docker Compose<sup>31</sup> is used to run the components listed in Table 1 as multi-container Docker applications. Docker Compose allows automatically retrieving or building required images and containers as well as controlling start up and shut down of multiple inter-dependent services. Docker Compose also allows the linking of services to make inter-service communication easier. Figure 5 shows the YAML file which defines

---

<sup>17</sup> <https://www.docker.com>

<sup>18</sup> <https://docs.docker.com/engine/installation>

<sup>19</sup> <http://www.mysql.com>

<sup>20</sup> Based on earkweb image created from <http://github.com/eark-project/earkweb>

<sup>21</sup> <https://lucene.apache.org/solr>

<sup>22</sup> [https://hub.docker.com/\\_/solr](https://hub.docker.com/_/solr)

<sup>23</sup> <http://www.rabbitmq.com>

<sup>24</sup> <https://hub.docker.com/r/tutum/rabbitmq>

<sup>25</sup> <http://redis.io>

<sup>26</sup> <https://github.com/tutumcloud/redis>

<sup>27</sup> <http://github.com/eark-project/earkweb>

<sup>28</sup> <https://github.com/eark-project/earkweb/blob/master/Dockerfile>

<sup>29</sup> <http://www.celeryproject.org>

<sup>30</sup> <https://github.com/mher/flower>

<sup>31</sup> <https://github.com/docker/compose>

the service composition so the services can be run together in an isolated environment. The dotted lines are used to visually divide the YAML file according to the different services it defines.

To give one example, the first service, named 'db', provides the MySQL database service needed by the *earkweb* frontend web application to store basic information about the processing status of information packages. The service is defined by the following properties:

- **image:** earkdbimg – The name of the image which is used to provide the MySQL database.
- **container\_name:** earkdb\_1 – The name of the Docker container.
- **build:** ./docker/earkdb – The Docker file which contains the build instructions to create the Docker image.
- **volumes:** - /tmp/earkweb-mysql-data:/var/lib/mysql – A directory from the Host system (here: /tmp/earkweb-mysql-data) is mounted as the MySQL data directory into the Docker container (here: /var/lib/mysql).
- **ports:** - '3306:3306' – The Port of the application is mapped to the port where the service will be exposed (in this case the standard MySQL port 3306 will be exposed as port 3306 by the container).

Furthermore, the orange arrows in Figure 5 show how the services are linked to each other using the "links" attribute of the *earkweb* container to reference the required service by name or the `BROKER_URL` environment variable of the flower service to link to the 'rabbitmq' service.

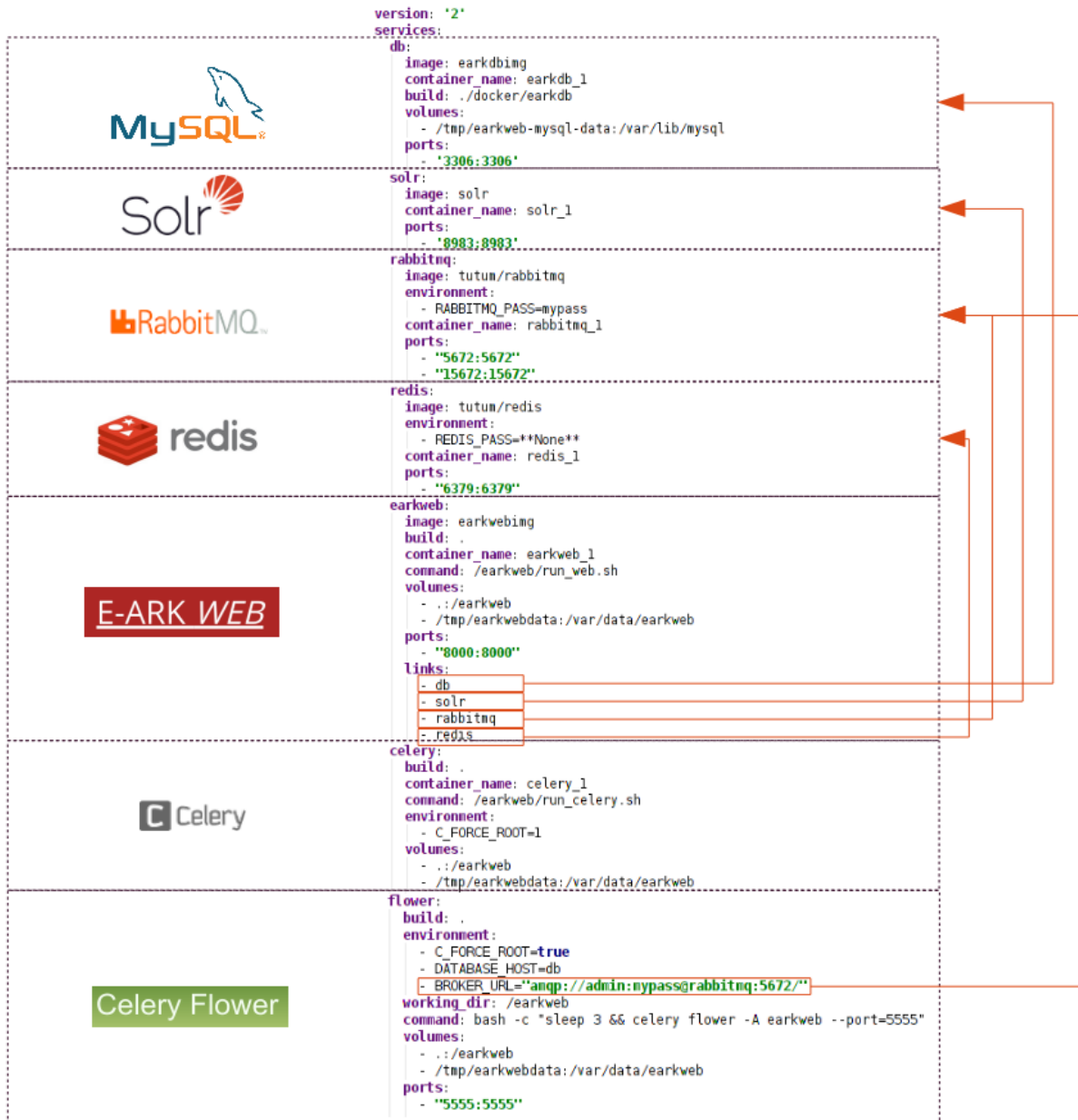


Figure 5 Docker containers of the Standalone Software Stack and the main links between the containers.

Figure 6 shows the origin of the images which are needed to run the various Docker containers. It shows that the Solr, RabbitMQ, and Redis images are directly retrieved from the Docker Image Library<sup>32</sup>, and that the

<sup>32</sup> <https://hub.docker.com>

other containers are based on Docker container build instructions provided in the form of Dockerfiles<sup>33</sup> as part of the *earkweb* code base.<sup>34</sup>

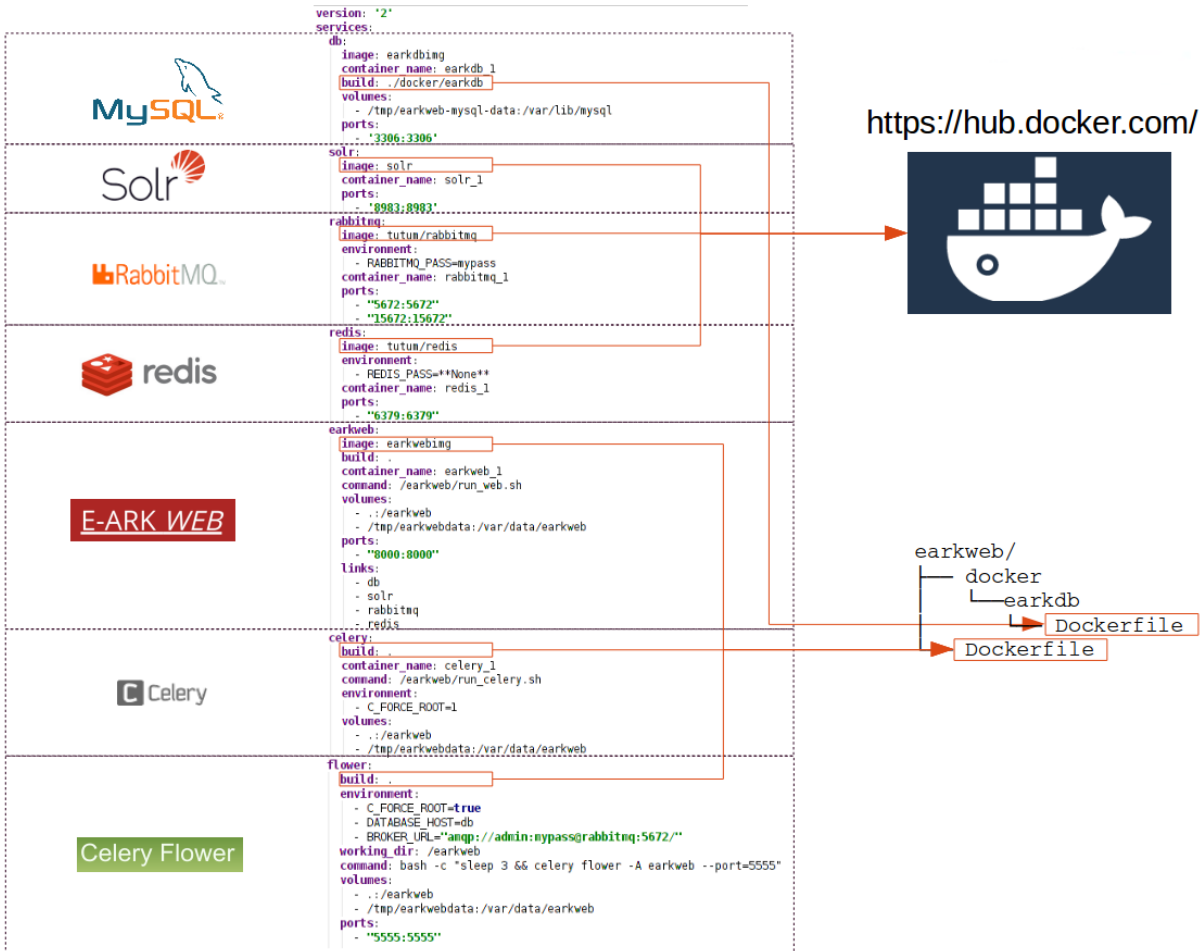


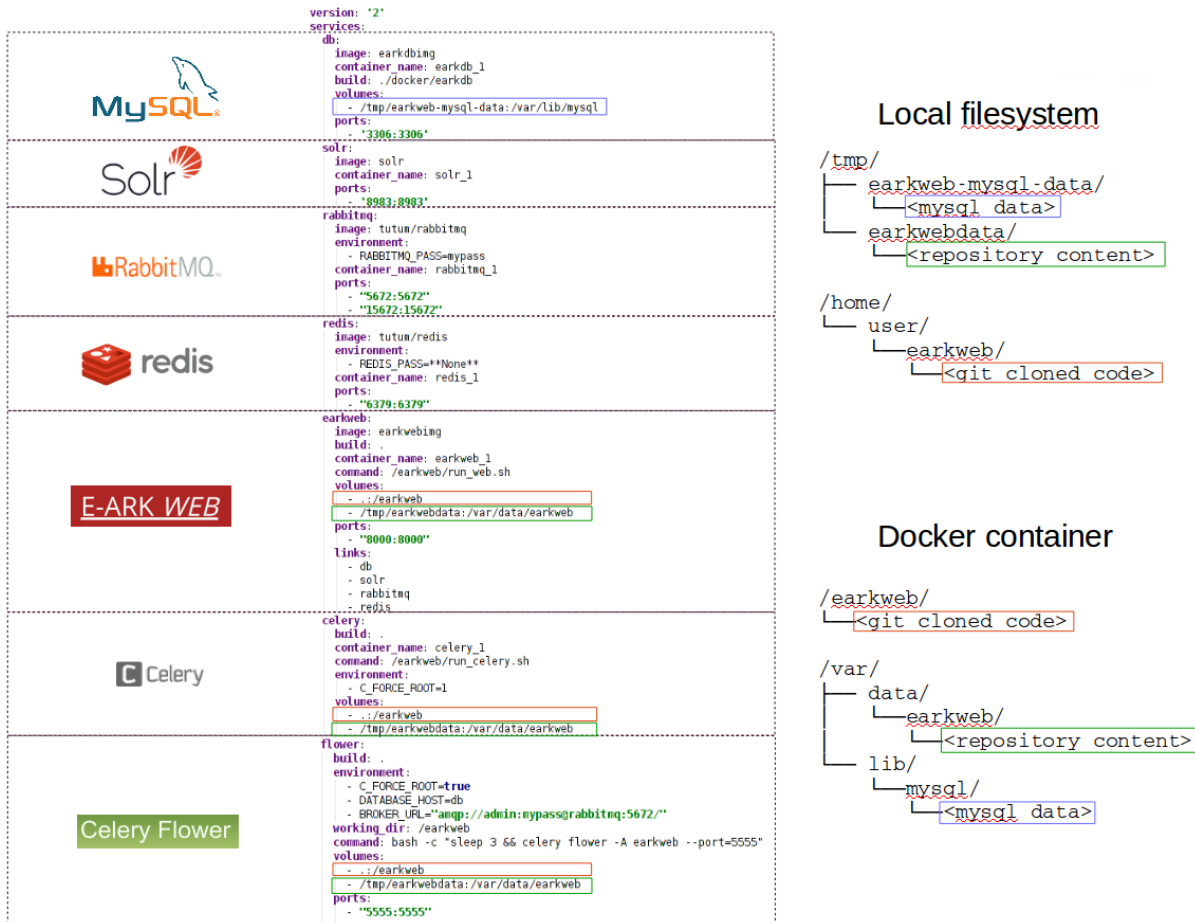
Figure 6 Origin of the images which are the basis to run the Docker containers

Finally, Figure 7 shows how various directories of the host file system are mapped to directories of the corresponding Docker containers. On the top right, there are directories of the local file system, and on the bottom right, there are directories of the Docker containers. The colours of the bounding boxes show the relationship between host file system and container directories. For example, the `/tmp/earkweb-mysql-data` directory of the host file system is mapped to the `/var/lib/mysql` directory of the MySQL database container. It is worth highlighting that for the containers *earkweb*, *Celery*, and *Celery Flower* the current directory (`.`) –

<sup>33</sup> <https://docs.docker.com/engine/reference/builder>

<sup>34</sup> <http://github.com/eark-project/earkweb>

which corresponds to the *earkweb* code base – is mounted to the */earkweb* container directory. These containers provide the different types of services (the *earkweb* frontend, the Celery task queue backend, and the Celery Flower task execution monitoring service) based on the same base image using service specific container run instructions respectively.



## 5 SIP-AIP conversion component

### 5.1 Task execution framework

The SIP-AIP conversion component consists of a set of individual tasks which are executed in a specific order to convert an E-ARK submission information package (SIP) into the E-ARK archival information package (AIP). It is an extensible workflow which can be adapted to specific needs by inserting new tasks at any point of the

workflow. *earkweb* uses a modular approach for defining atomic tasks which perform a specific transformation step of the SIP-AIP conversion, such as the extraction of an SIP or the validation of the descriptive metadata it contains. However, a specific task does not necessarily execute one single action, but can initiate a series of tasks or a complete workflow as well.

Each task is implemented as a python class and is available in the python module “workers/tasks.py”<sup>35</sup> of the *earkweb* application. A task which performs a step of the SIP-AIP conversion must extend the default task class `DefaultTask` defined in the module “workers/default\_task.py”.<sup>36</sup>

The default task makes sure that the pre-conditions for executing a task are fulfilled (e.g. the package is not in an error state), and also records the task execution in the provenance metadata of the package (i.e. the PREMIS metadata file). The default task also verifies if task execution is allowed given the current state of the package. Each task has a property which defines the list of tasks which are accepted as previously executed tasks. The fact that a task is defined as an “accepted last task” means that if execution was successful, there is the assumption that it produces valid output to be used as input for the current task. For example, to execute the `SIPValidation` task<sup>37</sup> which validates the structure of the E-ARK IP, it is a requirement that the transferred entity (the SIP’s TAR file) was extracted successfully by the `SIPExtraction` task.<sup>38</sup>

## 5.2 SIP-AIP conversion tasks

Table 2 provides an overview of the tasks which together represent the SIP-AIP conversion component.

| Task name                          | Accepted inputs  | Task description  |
|------------------------------------|--|---|
| <code>SIPDeliveryValidation</code> | <code>SIPValidation</code> ,<br><code>SIPtoAIPReset</code> | Validation of the delivery of the ZIP or TAR file against a delivery METS XML file to verify the integrity of the transferred file. |
| <code>IdentifierAssignment</code>  | <code>SIPDeliveryValidation</code>                         | Assignment of a unique identifier. In the reference   |

<sup>35</sup> <https://github.com/eark-project/earkweb/blob/master/workers/tasks.py>

<sup>36</sup> [https://github.com/eark-project/earkweb/blob/master/workers/default\\_task.py](https://github.com/eark-project/earkweb/blob/master/workers/default_task.py)

<sup>37</sup> <https://github.com/eark-project/earkweb/blob/6be44bc8ed3d346141c7f7e4091f1069c7a467f5/workers/tasks.py#L1100>

<sup>38</sup> <https://github.com/eark-project/earkweb/blob/6be44bc8ed3d346141c7f7e4091f1069c7a467f5/workers/tasks.py#L1059>



|                                  |  |   |
|----------------------------------|--|---|
|                                  |  | implementation a URN with the namespace identifier “uuid” is generated.   |
| SIPExtraction                    | IdentifierAssignment   | Extraction of the TAR or ZIP file containing the SIP.   |
| SIPValidation                    | SIPExtraction  | Validation of the structure of the SIP as well as the structural (METS) and preservation (PREMIS) metadata.   |
| SIPRestructuring                 | SIPValidation,<br>SIPExtraction  | Restructuring of the SIP according to the AIP structure which encloses the SIP.   |
| AIPDescriptiveMetadataValidation | AIPDescriptiveMetadataValidation,<br>SIPRestructuring  | Validation of descriptive metadata. In the current state, only EAD descriptive metadata is validated.   |
| AIPMigrations                    | SIPRestructuring,<br>AIPDescriptiveMetadataValidation,<br>MigrationProcess,<br>AIPMigrations,<br>AIPCheckMigrationProgress | Execute AIP migrations according to defined migration rules. As examples, PDF to PDF/A and GIF to TIFF migrations are defined with corresponding tools performing the migrations. |
| AIPCheckMigrationProgress        | AIPCheckMigrationProgress,<br>AIPMigrations,<br>MigrationProcess   | This task can be invoked to query the migration process.  |
| CreatePremisAfterMigration       | CreatePremisAfterMigration,<br>AIPCheckMigrationProgress   | Create PREMIS metadata after migration.   |
| AIPRepresentationMetsCreation    | AIPRepresentationMetsCreation,<br>CreatePremisAfterMigration   | Create the METS structural metadata file for each representation contained in the AIP.  |

|                        |   |   |
|------------------------|---|---|
| AIPPackageMetsCreation | AIPPackageMetsCreation,<br>AIPRepresentationMetsCreation,<br>CreatePremisAfterMigration | Create the AIP METS file referencing the METS files of the representations contained in the AIP.  |
| AIPValidation          | AIPValidation,<br>AIPPackageMetsCreation  | Validation of the structure of the AIP as well as the structural (METS) and preservation (PREMIS) metadata.   |
| AIPPackaging           | AIPPackaging,<br>AIPValidation  | Packaging the AIP as a TAR file.  |
| AIPStore               | AIPStore,<br>AIPPackaging   | Store the AIP in the file system in the Pairstree storage. This is the storage area of the standalone software stack. In the cluster software stack this storage area represents the staging area holding packages which are to be uploaded to the Lily Repository. <sup>39</sup> |
| LilyHDFSUpload         | AIPStore  | Upload AIP package file to the Lily Repository – only used in the cluster deployment stack.   |

Table 2 SIP-AIP conversion tasks

<sup>39</sup> The files are uploaded to the Hadoop Distributed File System (HDFS) of the Lily Repository. See deliverable D6.2 for details about the Lily Repository at <http://e-ark-project.com/resources/project-deliverables/54-d62intplatformref-1>

## 6 E-ARK Web user guide

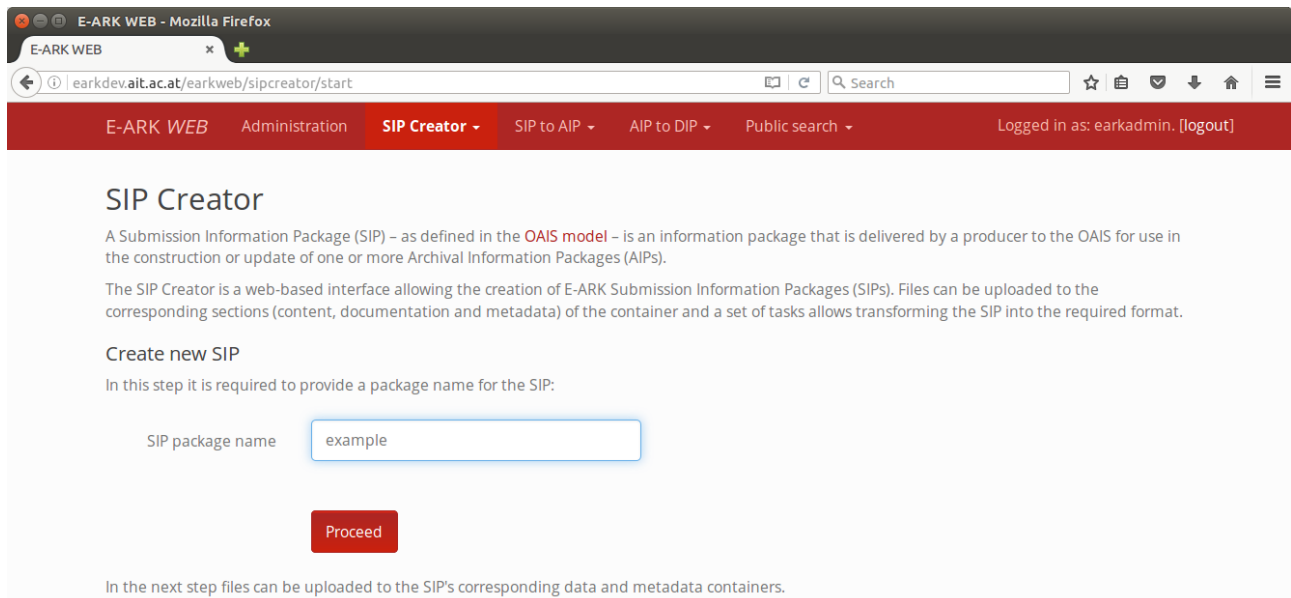
### 6.1 SIP Creator

A Submission Information Package (SIP) – as defined in the OAIS model – is an information package that is delivered by a producer to the OAIS for use in the construction or update of one or more Archival Information Packages (AIPs).

The SIP Creator is a web-based interface allowing the creation of E-ARK Submission Information Packages (SIPs). Files can be uploaded to the corresponding sections (content, documentation and metadata) of the container and a set of tasks allows transforming the SIP into the required format.

#### 6.1.1 Initialise SIP package creation

In this step it is required that a package name for the SIP be provided:



The screenshot shows a web browser window titled 'E-ARK WEB - Mozilla Firefox' with the address bar containing 'earkdevait.ac.at/earkweb/sipcreator/start'. The page has a red navigation bar with links for 'E-ARK WEB', 'Administration', 'SIP Creator', 'SIP to AIP', 'AIP to DIP', and 'Public search'. The user is logged in as 'earkadmin'. The main content area is titled 'SIP Creator' and contains the following text: 'A Submission Information Package (SIP) – as defined in the OAIS model – is an information package that is delivered by a producer to the OAIS for use in the construction or update of one or more Archival Information Packages (AIPs). The SIP Creator is a web-based interface allowing the creation of E-ARK Submission Information Packages (SIPs). Files can be uploaded to the corresponding sections (content, documentation and metadata) of the container and a set of tasks allows transforming the SIP into the required format.' Below this is a section titled 'Create new SIP' with the instruction 'In this step it is required to provide a package name for the SIP:'. There is a text input field labeled 'SIP package name' containing the text 'example'. Below the input field is a red button labeled 'Proceed'. At the bottom of the form, it says 'In the next step files can be uploaded to the SIP's corresponding data and metadata containers.'

Note that a package name must have at least 3 characters.

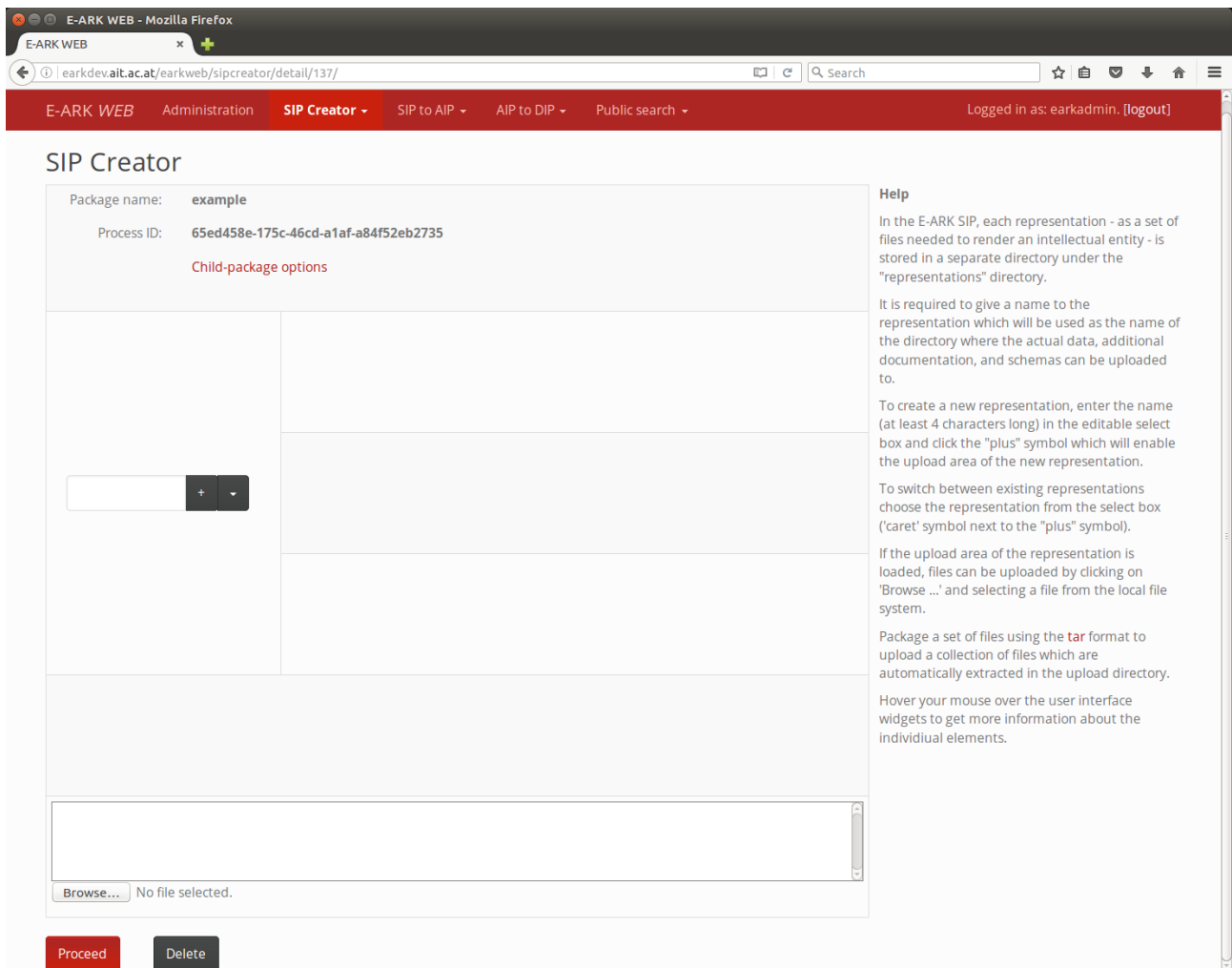
Click on 'Proceed' to continue with the next step.

#### 6.1.2 Adding files to the SIP

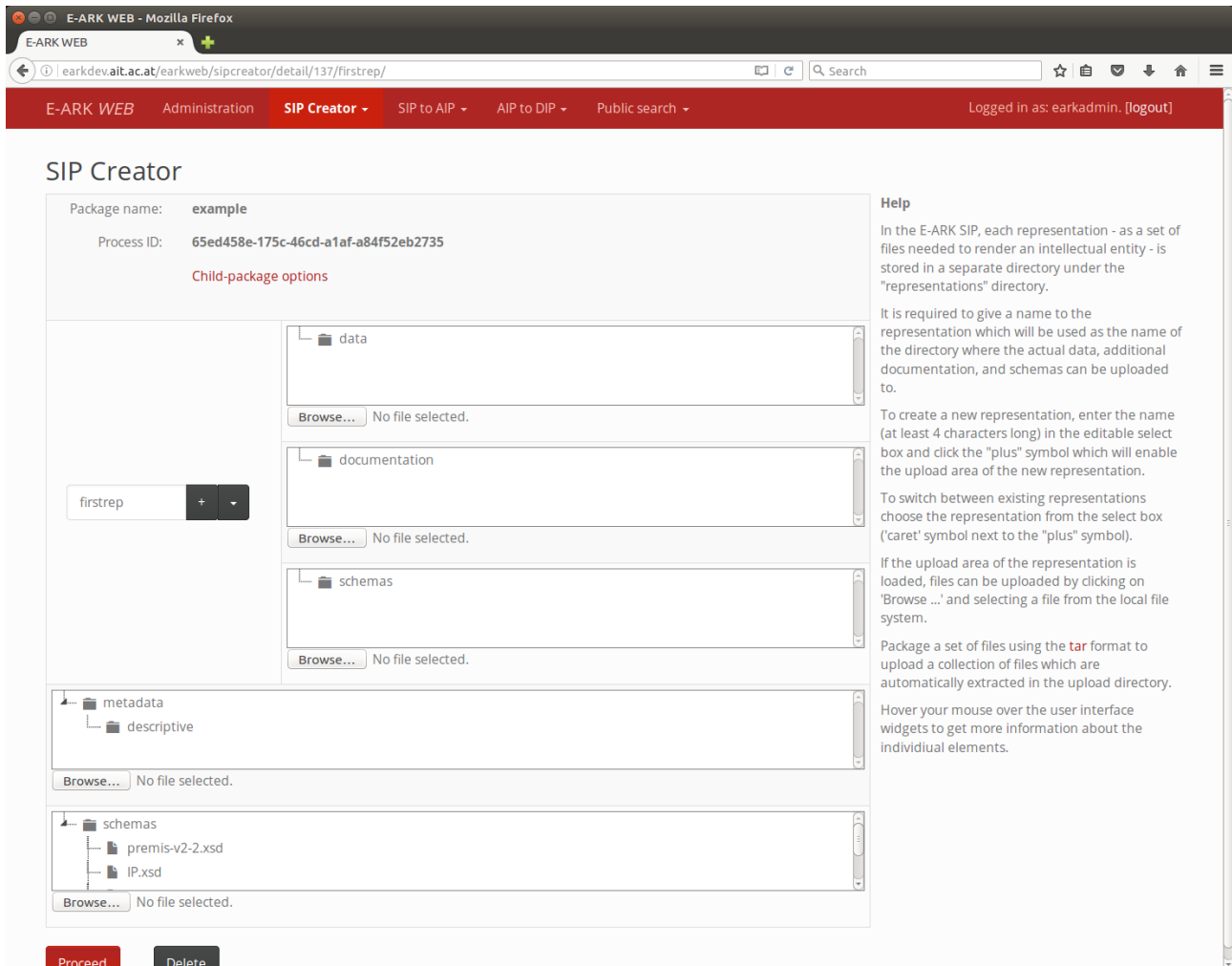
A web form allows uploading files to the corresponding areas of the SIP.

In the E-ARK SIP, each representation - as a set of files needed to render an intellectual entity - is stored in a separate directory under the "representations" directory.

It is required that a name be given to the representation which will be used as the name of the directory where the actual data, additional documentation, and schemas can be uploaded.



To create a new representation, enter the name (at least 4 characters long) in the editable select box and click the "plus" symbol which will enable the upload area of the new representation:



To switch between existing representations choose the representation from the select box ('caret' symbol next to the "plus" symbol).

If the upload area of the representation is loaded, files can be uploaded by clicking on 'Browse ...' and selecting a file from the local file system.

Package a set of files using the tar format to upload a collection of files which are automatically extracted in the upload directory.

Hover your mouse over the user interface widgets to get more information about the individual elements.

Click on the "Proceed" button to continue with the next step.

### 6.1.3 SIP creation process

#### 6.1.3.1 Start SIP creation process

The SIP creation process interface allows executing SIP creation tasks:

- SIPReset
- SIPDescriptiveMetadataValidation
- SIPPackageMetadataCreation
- SIPPackaging
- SIPClose

##### 6.1.3.1.1 *Task execution*

The tasks in the pull-down menu near the label "Tasks" appear in the logical order in which they should be executed.

**SIP Creator**

A Submission Information Package (SIP) – as defined in the **OAIS model** – is an information package that is delivered by a producer to the OAIS for use in the construction or update of one or more Archival Information Packages (AIPs).

The SIP Creator is a web-based interface allowing the creation of E-ARK Submission Information Packages (SIPs). Files can be uploaded to the corresponding sections (content, documentation and metadata) of the container and a set of tasks allows transforming the SIP into the required format.

**SIP creation task/workflow execution**

|                   |   |
|-------------------|---|
| Package name      | example   |
| Process ID        | 65ed458e-175c-46cd-a1af-a84f52eb2735  |
| Working area path | <a href="#">/var/data/earkweb/work/65ed458e-175c-46cd-a1af-a84f52eb2735</a> |
| Last task         | SIPReset  |
| Last change       | 09.05.2016 10:47:28   |
| Process status    | Success (0)   |

**Task/Workflow execution**

Task:

Process log

Error log

[back](#)

The "SIP creation task/workflow execution" overview table shows the information entities of the package:

- The "Package name" is the name of the SIP which is the outcome of the current SIP creation process.
- The "Process ID" is a UUID which corresponds to the name of the working area directory.
- The "Working area path" is a link which allows accessing the working area directory of the corresponding package to see the result of a task execution.
- The "Last task" shows the last task which was executed, e.g. "SIPValidation". The last task also determines the following task which can be executed because each task defines a list of accepted input tasks.
- Last change is the date/time of the last modification (by a task)
- Process status shows the consistency status of a package. If the process status shows the value

"Success (0)" together with a green "check mark" symbol, then the status of the package is consistent. If an error occurred during processing, the process status shows the value "Error (1)" together with a green "warning triangle" symbol.

Selected information and error log messages appear in the "Process log" and "Error log" areas, more detailed information about the processing might be available in the *earkweb* processing log which is available in the package.

Use the "SIPReset" task to roll-back package processing to the initial state in case an error occurred during processing (after applying required modifications).

It is possible to directly continue with the "SIP to AIP conversion" process by executing the "SIP Close" task. Note that in this case the package disappears from the "Active SIP creation processes" overview page as it is handed over to the "SIP to AIP conversion".

#### 6.1.3.2 Active SIP creation processes overview

By accessing the sub-menu item "Active SIP creation processes", an overview about open SIP creation processes is shown. Each package is identified by the package name (column 'Package name') which was provided in the first step of the SIP creation process and an internal process identifier (column 'Process ID'). The process identifier is also used as the name of the working directory where information package transformations take place (column 'SIP creation working directory').



**SIP Creator**

A Submission Information Package (SIP) – as defined in the **OAIS model** – is an information package that is delivered by a producer to the OAIS for use in the construction or update of one or more Archival Information Packages (AIPs).

The SIP Creator is a web-based interface allowing the creation of E-ARK Submission Information Packages (SIPs). Files can be uploaded to the corresponding sections (content, documentation and metadata) of the container and a set of tasks allows transforming the SIP into the required format.

**Active SIP creation processes**

The following table gives an overview about open SIP creation processes. Each package is identified by the package name (column 'Package name') which was provided in the first step of the SIP creation process and an internal process identifier (column 'Process ID'). The process identifier is also used as the name of the working directory where information package transformations take place (column 'SIP creation working directory').

Depending on the last task that was executed, subsequent transformation tasks can be applied.

| SIP package name     | Process ID                           | Last task                  | Last change         | Status  |
|----------------------|--------------------------------------|----------------------------|---------------------|---------|
| example              | 65ed458e-175c-46cd-a1af-a84f52eb2735 | SIPReset                   | 09.05.2016 10:47:28 | Success |
| ERMS2                | 66d8d25b-3cee-4475-9cc9-f69bfd7e64a  | SIPPackaging               | 27.04.2016 09:59:09 | Success |
| zip                  | 69792876-2b87-4208-9e22-3b1d27699fbc | SIPReset                   | 22.04.2016 20:56:04 | Success |
| ID.AVID.18005        | 7fbed40c-5a55-4f58-b785-34b90fdefab  | SIPReset                   | 01.04.2016 11:46:43 | Success |
| CAMBRIDGE.001        | d5dd80d6-4673-4de4-80da-977b7d7c0e35 | SIPReset                   | 12.03.2016 11:28:47 | Success |
| ID.AVID.RA.18005.001 | d40bab5e-b057-4985-af2d-ee798ccba3d8 | SIPPackaging               | 08.03.2016 16:33:55 | Success |
| Paket.123            | 7f8cd77-e6c2-4b65-bb39-907c2d17ae59  | SIPPackaging               | 21.01.2016 12:47:09 | Success |
| mnl.pack.001         | aaac9706-3969-4ac3-8129-66e43b76968a | SIPReset                   | 13.01.2016 11:23:25 | Success |
| friendsinsweden      | 8cb0d227-ed0a-4376-91b7-9bda1979e1af | SIPReset                   | 15.12.2015 15:36:18 | Success |
| PARENT               | 04838029-0d59-4300-bd32-906d3b0b3e9a | SIPReset                   | 14.12.2015 18:33:10 | Success |
| P.SIP                | d59d77a1-ddc5-4431-80ac-8f4f6e2625ae | SIPPackaging               | 03.12.2015 14:07:48 | Success |
| ERA.5.4.3            | 49802e81-45a2-45bf-87d9-dd358242dccb | SIPReset                   | 02.12.2015 14:21:01 | Success |
| EARK.SIP.001         | e1e0126d-bd50-4d76-a759-eb8f3e237322 | SIPPackaging               | 24.11.2015 12:55:40 | Success |
| PREMISTEST           | 2a519d4b-618b-4a25-99e8-8a5ab25de8fe | SIPPackageMetadataCreation | 19.11.2015 16:15:04 | Success |
| xpto                 | 9e9a1040-358d-4fc6-b8b9-90c763fec55c | SIPReset                   | 19.11.2015 15:29:28 | Success |
| Beemen_test_2        | 7724845b-6980-4acc-8804-6c6b015b6fee | SIPReset                   | 05.11.2015 10:17:18 | Success |

Depending on the last task that was executed, subsequent transformation tasks can be applied.

## 6.2 SIP to AIP conversion

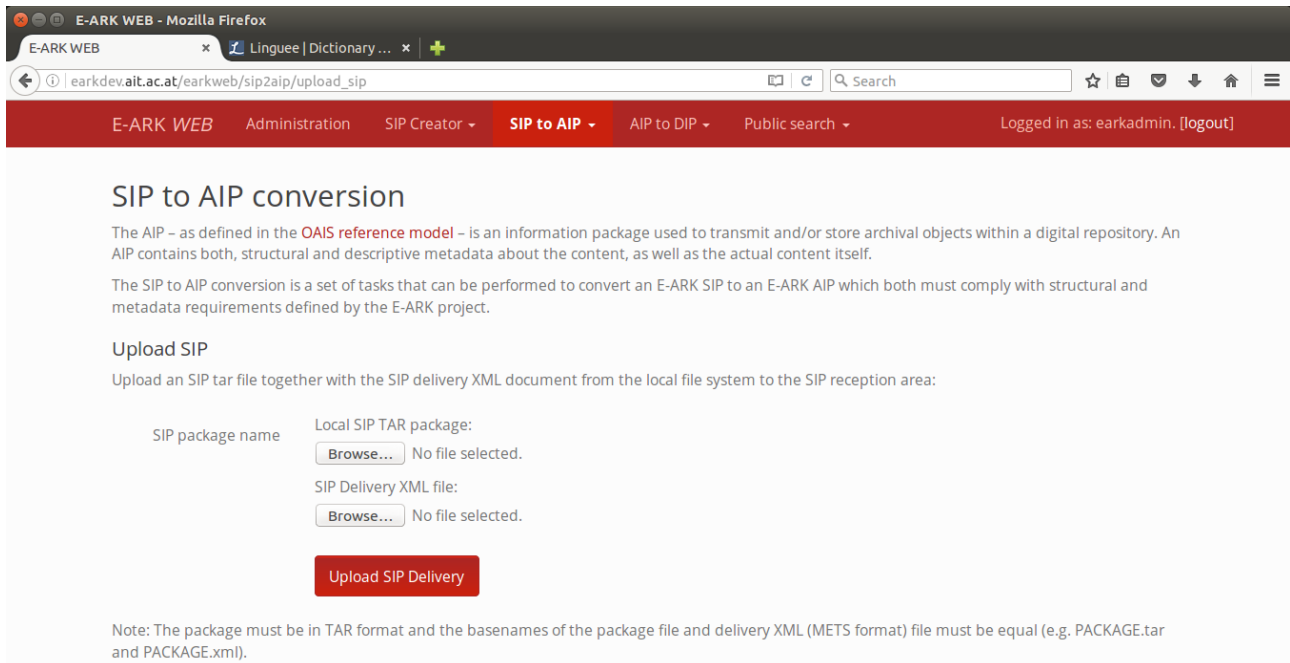
The AIP – as defined in the OAIS reference model – is an information package used to transmit and/or store archival objects within a digital repository. An AIP contains both, structural and descriptive metadata about the content, as well as the actual content itself.

The SIP to AIP conversion is a set of tasks that can be performed to convert an E-ARK SIP to an E-ARK AIP which must comply with both structural and metadata requirements defined by the E-ARK project.

### 6.2.1 Start SIP to AIP conversion process

The "SIP to AIP conversion" is either started by handing the package over from the "SIP creation" process

within *earkweb* by executing the "SIP Close" task or it can be uploaded using the "SIP to AIP conversion" sub-menu item "Upload SIP":



The TAR or ZIP container file of an SIP is uploaded together with the SIP delivery XML document from the local file system to the SIP reception area.

Note that the package must be in TAR or ZIP format and the basenames of the package file and delivery XML (METS format) file must be equal (e.g. PACKAGE.tar and PACKAGE.xml).

### 6.2.1.1 SIP to AIP conversion tasks

The SIP to AIP conversion process interface allows executing SIP creation tasks:

- SIPToAIPReset - Rollback to initial state (as after handover from SIP creation or uploading the SIP)
- SIPDeliveryValidation - Validate delivery (especially checksum information)
- IdentifierAssignment - Assign identifier
- SIPExtraction - Extract SIP
- SIPRestructuring - Restructure SIP according to AIP format
- SIPValidation - Validate SIP
- AIPMigrations - Do migrations according to predefined rules (e.g. PDF to PDF/A)

- AIPCheckMigrationProgress - Check if migrations are finished
- CreatePremisAfterMigration - Create PREMIS file after migrations completed
- AIPRepresentationMetsCreation - Create representation METS files
- AIPPackageMetsCreation - Create AIP METS and PREMIS
- AIPValidation - Validate package
- AIPPackaging - Create package container file
- AIPStore - Store AIP package in file system storage area (pairtree storage)
- AIPIndexing - Index AIP
- LilyHDFSUpload - Upload AIP to public search server

#### 6.2.1.2 SIP to AIP conversion task execution

The tasks in the pull-down menu near the label "Tasks" appear in the logical order in which they should be executed.

The "SIP to AIP task/workflow execution" overview table shows the information entities of the package. Note that some information entities, such as "identifier" might appear only after a specific task ("IdentifierAssignment" in this case), was executed.

- The "Process ID" is a UUID which corresponds to the name of the working area directory.
- The "Package name" is the name of the SIP which is the outcome of the current SIP creation process.
- The package Identifier is the PUID of the AIP. Although of type UUID in the reference implementation, this could be any other type of identifier, such as a DOI or PURL.
- The "Working area path" is a link which allows accessing the working area directory of the corresponding package to see the result of a task execution.
- The "Last task" shows the last task which was executed, e.g. "AIPValidation". The last task also determines the subsequent task which can be executed because each task defines a list of accepted input tasks.
- Last change is the date/time of the last modification (by a task).
- Process status shows the consistency status of a package. If the process status shows the value "Success (0)" together with a green "check mark" symbol, then the status of the package is consistent.

If an error occurred during processing, the process status shows the value "Error (1)" together with a green "warning triangle" symbol.

Selected information and error log messages appear in the "Process log" and "Error log" areas, more detailed information about the processing might be available in the *arkweb* processing log which is available in the package.

Use the "SIPtoAIPReset" task to roll-back package processing to the initial state in case an error occurred during processing (after applying required modifications).

### 6.2.1.3 Active SIP to AIP conversion processes

By accessing the sub-menu item "Active SIP to AIP conversion processes", an overview about open SIP to AIP conversion processes is shown. Each package is identified by the package name (column 'Package name') which corresponds to the name of the SIP container. The process identifier is also used as the name of the working directory where information package transformations take place (column 'SIP creation working directory').

**SIP to AIP conversion**

The AIP – as defined in the [OAI reference model](#) – is an information package used to transmit and/or store archival objects within a digital repository. An AIP contains both, structural and descriptive metadata about the content, as well as the actual content itself.

The SIP to AIP conversion is a set of tasks that can be performed to convert an E-ARK SIP to an E-ARK AIP which both must comply with structural and metadata requirements defined by the E-ARK project.

**Active SIP to AIP conversion processes**

The following table gives an overview about open SIP to AIP conversion processes. Each package has a package name which was provided in the first step of the SIP creation process and an internal process identifier (Process ID). The process identifier is also the name of the working directory where information package transformations take place (Working directory).

| Package name      | Process ID                           | Last task              | Last change         | Process status |
|-------------------|--------------------------------------|------------------------|---------------------|----------------|
| EINSTEINSTATE.001 | 5c23ac70-787f-4f58-9a08-b2cd29a35936 | AIPPackaging           | 29.04.2016 21:18:32 | Success        |
| TESTSIP           | b938322b-69e8-47fb-9e0a-ae39c2f27f99 | SIPValidation          | 27.04.2016 10:35:01 | Error          |
| METS              | 5308f797-57ce-45ad-ab93-3ac40f238add | SIPDeliveryValidation  | 27.04.2016 10:02:50 | Success        |
| ERMS.SMURF.1      | f1e4f558-44eb-4456-8168-1a041a705d69 | SIPClose               | 20.04.2016 09:44:53 | Success        |
| REVIEW            | 1e126d6c-209d-49a6-9bd6-135803080c02 | LilyHDFSUpload         | 14.04.2016 15:27:29 | Success        |
| EINSTEINAIP.001   | ad96fd64-9418-40eb-b8a2-e60f1a0eaa9a | LilyHDFSUpload         | 13.04.2016 15:58:20 | Success        |
| aaaa              | 83d87e81-aa11-444a-a500-4f8fe651a69e | LilyHDFSUpload         | 12.04.2016 13:31:04 | Success        |
| VAUBAN.001        | 405fe771-dbfd-4a4c-965a-99ef963aba30 | SIPValidation          | 21.03.2016 14:09:35 | Error          |
| Madrid-01         | 1c8ee501-b156-4b99-80bb-511c1b3ed6be | LilyHDFSUpload         | 01.03.2016 18:03:29 | Success        |
| politikpdf        | 7e4a1162-bc9d-4d85-87b5-29d9ae258ce6 | LilyHDFSUpload         | 15.02.2016 17:22:54 | Success        |
| pdfsforcif        | b84b478e-71c2-4c82-b701-59b3b8736ca3 | LilyHDFSUpload         | 15.02.2016 16:53:58 | Success        |
| techmeeting       | f0f73a13-7c43-4bdd-81a7-bb3ce3ed3a79 | AIPPackageMetsCreation | 14.12.2015 13:53:33 | Success        |
| SIARD3            | 31463462-dab4-4a8e-b9be-d27937189231 | AIPPackaging           | 11.12.2015 10:29:23 | Success        |
| SLOV.GEO.ALL      | 5c6f5563-7665-4719-a2b6-4356ea033c1d | SIPRestructuring       | 10.12.2015 17:17:03 | Success        |

Depending on the last task that was executed, subsequent transformation tasks can be applied.

#### 6.2.1.4 Indexing and search in AIPs

The "AIPIndexing" task allows indexing the AIP TAR package created the SIP to AIP conversion process. The AIP is indexed according to the location in the pairtree storage. In that way the search results can offer a link to directly render individual content items.

### 6.2.1.5 AIPs in pairtree storage

The storage backend used in the Standalone Software Stack<sup>40</sup> of *earkweb* is based on the Python-based *Pairtree File System* implementation<sup>41</sup> of the *Pairtrees for Collection Storage* specification<sup>42</sup>. It is used to store the physical representation of the AIP in a conventional file system as opposed to the Hadoop Distributed File System which is used in the Revised Cluster Software Stack described in Deliverable D6.3.

In summary, the *Pairtree* is a filesystem hierarchy for storing digital objects where the unique identifier string of the object is mapped to a unique directory path so that the file system location of the object can be derived from the identifier string. Basically, the identifier string is split each two characters at a time and the object folder has by definition more than two characters. Furthermore, the specification defines a mapping of special characters to a set of alternative characters in order to ensure file system level interoperability.

The following example will explain how the *Pairtree* storage method is used in the *earkweb* implementation.

According to the default *earkweb* configuration, the path to the storage directory in the file system is:

```
/var/data/earkweb/storage
```

The storage folder contains an empty file called “pairtree\_version0\_1” which specifies the version of the *Pairtree* specification.

The storage directory contains the root folder of the *Pairtree* storage:

```
/var/data/earkweb/storage/pairtree_root
```

Let us assume that an AIP has the following identifier

```
urn:uuid:6c496473-4e77-44f5-b387-25bffd362789
```

Following the method defined by the *Pairtree* specification, this identifier is mapped to the following path:

```
ur/n+/uu/id/+6/c4/96/47/3-/4e/77/-4/4f/5-/b3/87/-2/5b/ff/d3/62/78/9
```

and the actual AIP container file is stored in a “data” folder which represents the leaf node of the *Pairtree* file system hierarchy.

The leaf node contains one or possibly more sub-directories (5-digits fixed length zero-filled number) for the versions of the AIP.

---

<sup>40</sup> In contrast to the cluster software stack, the standalone software stack does not include a big data backend based on Apache Hadoop.

<sup>41</sup> <https://pypi.python.org/pypi/Pairtree>

<sup>42</sup> <https://wiki.ucop.edu/display/Curation/PairTree?preview=/14254128/16973838/PairtreeSpec.pdf>

The full path to the AIP file is as follows (to be read as a single line string):

```
/var/data/earkweb/storage/pairtree_root/ur/n+/uu/id/+6/c4/96/47/3-/4e/77/-  
4/4f/5-/b3/87/-2/5b/ff/d3/62/78/9/data/00001/urn+uuid  
+6c496473-4e77-44f5-b387-25bffd362789.tar
```

Note that the colon characters (":") which are part of the URN identifier to separate the leading "urn" sequence from the namespace identifier ("uuid") and the latter from the namespace-specific string ("6c496473-4e77-44f5-b387-25bffd362789") were mapped to the plus character ("+"). This is because the *Pairtree* character mapping rules apply to the filename of the physical container as well to ensure file system interoperability.

The purpose of the *Pairtree* storage backend is to allow a large number of AIPs<sup>43</sup> to be stored in a conventional file system, such as Network Attached Storage, for example, and allow fast access to individual files by directly reading content streams from the files stored in the TAR packaged AIP container files.

If the AIP is changed, because, for example, a new representation was added to the AIP, a new physical container will be created as a new version of the AIP. As a result, only the numerical folder will be incremented (00002) and the identifier for the intellectual entity remains the same (426087e8-0f79-11e3-847a-34e6d700c47b):

```
/var/data/earkweb/storage/pairtree_root/ur/n+/uu/id/+6/c4/96/47/3-/4e/77/-  
4/4f/5-/b3/87/-2/5b/ff/d3/62/78/9/data/00002/urn+uuid  
+6c496473-4e77-44f5-b387-25bffd362789.tar
```

#### 6.2.1.6 Parent-child AIP relations

During AIP creation, it is possible to assign parent-child relationships between AIPs, as described in D4.4 part A. Figure 8 illustrates how the creation of packages governed by a parent-child relationship can be achieved:

1. The submission of a package is processed up to the point where the identifier of the parent AIP ("parentID") is created (operation "submitParentSIP").
2. This SIP to AIP creation process remains open until the last child AIP is created.

---

<sup>43</sup> No concrete numbers will be given here about what a "large number of AIPs" means in this context, as the requirements regarding scalability of indexing and access procedures depend on various factors, such as the environment (hardware and network) and the type of collection data to be stored in the repository. It is a matter of evaluating the Standalone Software Stack (cf. section **Error! Reference source not found.**) using scalability tests to find out if it can cope with the given scalability requirements, or if the Cluster Software Stack (cf. section **Error! Reference source not found.**) is needed.

3. A sequence of child SIPs can then be submitted and created as child-AIPs by indicating the parent identifier (“parentID”) generated previously.
4. After the last child-AIP is processed, the parent-AIP is created (operation “createParentAIP”) by indicating which child-AIPs belong to the AIP (“childs”).
5. Finally, the consistency of the logical AIP composition is verified and the parent-AIP (“parent”) and all child-AIPs (“childs”) are stored (operation “storeAIP”).

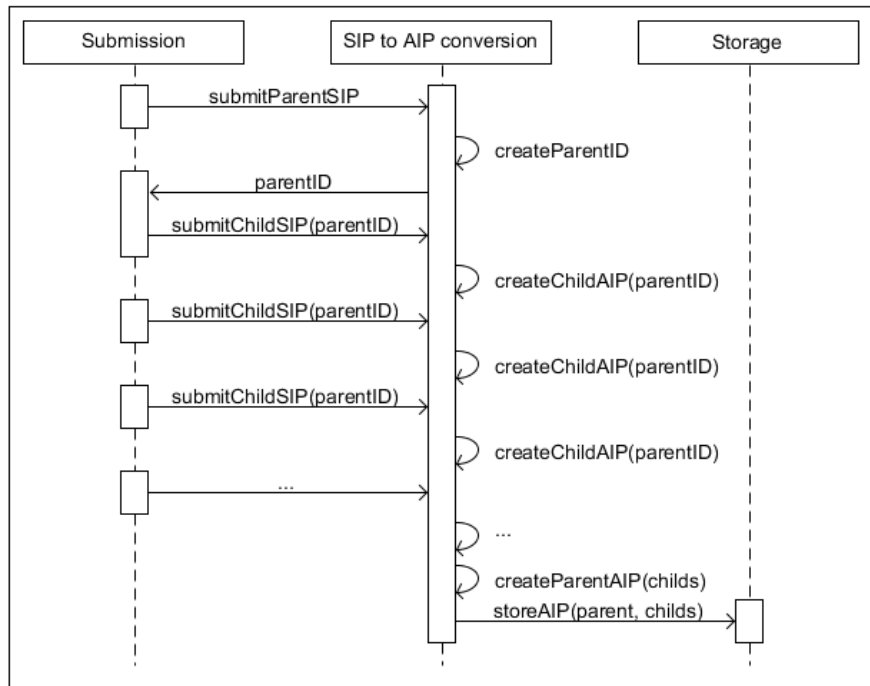


Figure 8: Parent-child aggregation submission process