

Deutscher Titel

DIPLOMARBEIT

zur Erlangung des akademischen Grades

Diplom-Ingenieur

im Rahmen des Studiums

Software Engineering and Internet Computing

eingereicht von

Bernhard Gößwein

Matrikelnummer 01026884

an der Fakultät für Informatik

der Technischen Universität Wien

Betreuung: Ao.Univ.Prof. Dipl.-Ing. Dr.techn. Andreas Rauber

Mitwirkung: Dr. Tomasz Miksa

Wien, 1. Jänner 2001

Bernhard Gößwein

Andreas Rauber

Designing a Framework gaining Repeatability for the OpenEO platform

DIPLOMA THESIS

submitted in partial fulfillment of the requirements for the degree of

Diplom-Ingenieur

in

Software Engineering and Internet Computing

by

Bernhard Gößwein

Registration Number 01026884

to the Faculty of Informatics

at the TU Wien

Advisor: Ao.Univ.Prof. Dipl.-Ing. Dr.techn. Andreas Rauber

Assistance: Dr. Tomasz Miksa

Vienna, 1st January, 2001

Bernhard Gößwein

Andreas Rauber

Erklärung zur Verfassung der Arbeit

Bernhard Gößwein
Vorderer Ödhof 1, 3062 Kirchstetten

Hiermit erkläre ich, dass ich diese Arbeit selbständig verfasst habe, dass ich die verwendeten Quellen und Hilfsmittel vollständig angegeben habe und dass ich die Stellen der Arbeit – einschließlich Tabellen, Karten und Abbildungen –, die anderen Werken oder dem Internet im Wortlaut oder dem Sinn nach entnommen sind, auf jeden Fall unter Angabe der Quelle als Entlehnung kenntlich gemacht habe.

Wien, 1. Jänner 2001

Bernhard Gößwein

Danksagung

Ihr Text hier.

Acknowledgements

Enter your text here.

Kurzfassung

Ihr Text hier.

Abstract

Enter your text here.

Contents

Kurzfassung	xi
Abstract	xiii
Contents	xv
1 Introduction	1
1.1 Problem Description	1
1.2 Aim of this Work	2
2 Related Work	3
2.1 Reproducibility	3
2.2 Related Work in the EO Sciences	4
2.3 Related Tools	4
2.4 OpenEO Details	4
3 Design	9
3.1 Back end provenance	10
3.2 Job dependent provenance	11
3.3 User Information	13
4 Implementation	15
4.1 The EODC Backend	15
4.2 Job Independent Context	15
4.3 Job Dependent Context	15
4.4 User Interface	15
4.5 Overview of Changes	15
5 Evaluation	17
5.1 Use Cases	17
5.2 Using Python Client	17
6 Conclusion and future Work	19
	xv

List of Figures	21
List of Tables	23
Index	25
Glossary	27
Acronyms	29
Bibliography	31

Introduction

1.1 Problem Description

Over the last decades remote sensing agencies have increased the variations of data processing and therefore the amount of resulting data. To preserve the data for further usage in the future it is necessary to have citable data and processes on the data to ensure repeatability in a long-term.[RMD12] Already most of the the data used in earth observation sciences are retrieved or provided via Service Oriented Architecture (SOA) interfaces. Provider like Google Earth Engine and EODC provide an Web API for retrieving and processing data. Due to a different range of functionality and a difference between the endpoints of the providers it is hard to create a workflow for more than one provider. The OpenEO project has the goal to be an abstraction layer above different EO data providers. The underlying structure consists of three parts:

- Client Module : Is written in the program language of the user and transfers the users commands to the backends.
- Core Module: A standard on how the communication should take place between client and backend.
- Backend Module: The provider of the data and the services, which gets the instructions from the clients and returns the results.

Further information on the software architecture of the project is defined in the project proposal ([PWS⁺17]).Until now there is no consideration of repeatability verification of workflows for users in the OpenEO architecture.Generalised layers have the opportunity to be implemented in a way that makes processes and data scientifically verifiable and reproducible, because it handles data and processes on the data in a standardised way on different providers. Even though the range of functionality and the API endpoints

are well defined in the OpenEO coreAPI the contributing content providers (OpenEO backends) will have different underlying software types and versions. The underlying technology of an OpenEO backend will also change over time and can lead to different results on the same workflow executions. Consider the following: A scientist runs an experiment using OpenEO as his research tool and gets results. The same scientist runs the same experiment with the same input some months later and gets slightly different results. The question occurs, why are the results different? Has the used data changed, has the user accidentally submitted different code or has some underlying software inside the backend provider changed. Adding a possibility for the users of OpenEO to gain this information is an important feature for the scientific community. The aim of this thesis is to provide a possibility for users of OpenEO to verify and validate a job re-execution on different underlying technologies of an OpenEO backend provider.[PWS⁺17]

Read through and improve above text.

1.2 Aim of this Work

The expected outcome of this thesis is to discover and develop a possible framework for providing repeatability in the OpenEO project. This enables users to re-execute workflows and validate the results, so that differences on the process or data are accessible for the users. To achieve this goal a model for repeatability within the project has to be discovered and implemented to evaluate the ability of the model. The model shall then conclude recommendations for the OpenEO project on how to improve re-execution validation for the user and how it can be achieved. Therefore the following research questions can be formulated:

- **How can an OpenEO job re-executed be applied like the initial execution?**
 - How can the used data be identified after the initial execution?
 - How can the used software of the initial execution be reproduced?
 - What data has to be captured when?
 - How can the result of a re-execution in future software versions be verified?
- **How can the equality of the OpenEO job re-execution results be validated?**
 - What are the validation requirements?
 - How can the data be compared?
 - How can the re-execution be validated after changes of the OpenEO backend environment?
 - How can differences in the environment between the executions be discovered?

Read through above text and improve. 2

Add description of use cases (see summary).

Related Work

Currently, there exists no concrete solution to add the ability of repeatability to the OpenEO project. However there are concepts of adding repeatability in computer science.

2.1 Reproducibility

2.1.1 eScience

The eScience has the potential to enable a boost in scientific discovery by providing approaches to make digital data and workflows citable. In [RMMP15] is a common way of reaching this goal formulated. It describes an approach to look at whole research processes, other than only data citation by introducing Process Management Plans. The capturing, verification and validation of the needed data for a computational process is also demonstrated within the paper.[RMMP15]

2.1.2 Data Citation

Since the earth observation community use a high amount of satellite data and also within the OpenEO project a lot of big data sets are being used, there needs to be a solution to cite the used data in a workflow. The Research Data Alliance (RDA) working group on data citation provides a 14 step recommendation of data citation. It contains solutions not only for static, but also for dynamic data, so data that changes over time. Using the guideline for data citations from the RDA makes the data scientifically citable. [RAvUP16] In earth science there is also a strategy of ESA and NASA to achieve a content standard for data preservation.[RMD12]

2.1.3 Provenance Data

The re-execution of an OpenEO workflow not only needs data citation, but also the information of how the workflow was executed. Therefore provenance data has to be captured.[RMH⁺11] There are already several provenance models defined in the scientific community. One of the existing models is the PROV model, which was published in 2013 by the World Wide Web Consortium Provenance Working Group and consists of recommendations and guidelines for provenance data.[MGC⁺15] Another model is the VFramework, designed for the purpose of redeployment including the verification of a re-execution of the same workflow. [MPM⁺13]

2.2 Related Work in the EO Sciences

2.3 Related Tools

Read through
above text and
improve.

2.3.1 ReproZip

Add section about
ReproZip.

2.3.2 Smartcontainer

Add section about
Smartcontainer.

2.3.3 Reproducibility in computational geoscience

2.4 OpenEO Details

The OpenEO project contains of three modules, the client modules written in the programming language of the users, the back end drivers that makes it possible for every backend to understand the calls from the clients and the core API that specifies how the communication should take place. So the core API is a standard that the back end providers accepted to implement on their systems. The back end drivers are the translation of the client calls to the back end specific API. This architecture decouples the clients from the back ends so that every client can connect to every back end that applies to the OpenEO core API standard see figure

The communication is specified as an OpenAPI(TODO: Quelle) description and consists of all RESTful services that can be called at the back ends. Most of the calls are only to receive information about the processes and the data of the specific back end, but there are also endpoints to apply job executions specified by the user. In this thesis the focus lies on the job execution, because this is the core feature of the OpenEO project and needs to be reproducible. In the next chapter the job execution process is described in more detail.

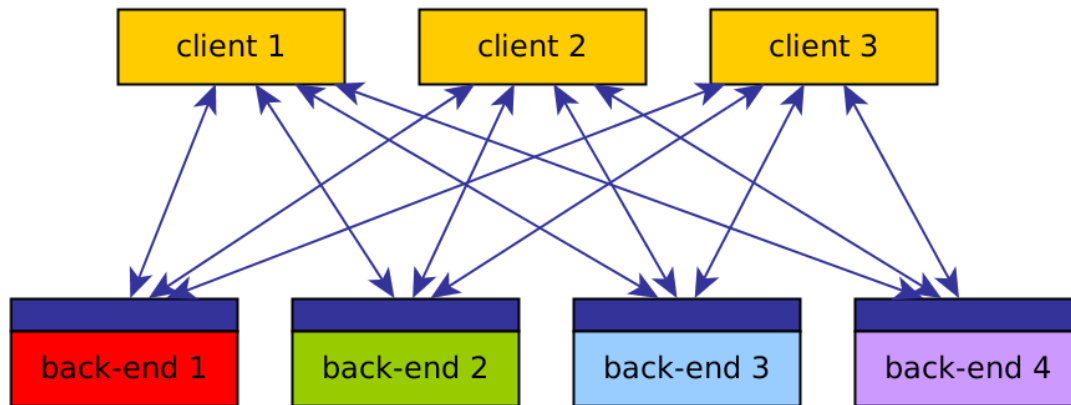


Figure 2.1: Overview of the OpenEO architecture. (from openeo.org)

2.4.1 Job Execution

The job execution workflow in OpenEO starts at a client application that let the user define what has to be processed in the specific programming language. The main part of the job execution is based on the description of what the back end needs to do with what data. Therefore OpenEO introduces the process graph, which consists of a tree structure describing the processes with their data and the input data identifier, which is back end specific. The process graph is in a JSON format and gets generated by the clients in the background without the users seeing it directly. In figure

The back ends interpret the process graph from inside out. Inside figure TODO the element in the middle defines the input data identifier in the `imagery` block, with the `product_id`. In this case the `s2a_prd_ms11c` is the chosen input data identifier. Now the back end goes one step up in the hierarchy of the process graph and calls the process `filter_bbox` with the parameters `left`, `right` and so on using the input data. The output data of the previous process is also the input data of the next process. It is possible that the input of a process is a whole new process graph. In figure

Add Diagramm of
Process Graph

There are two different kinds of process executions depending on the back ends capabilities, synchronous and asynchronous calls. Synchronous calls are directly executed after they got received from the back end and the user that sends the process graph has to wait until the job is finished. For example on the python client the program waits after sending the process graph to the back end until the back end returns the result and the results are directly returned to the user. On the asynchronous call the user sends the process graph, but it does not get executed until the user starts the execution on the back end through an additional endpoint. When the processing is finished the user can download the result at another endpoint of the back end. For the asynchronous calls there is also the possibility to subscribe to a notification system on the backend, so that the user gets

2. RELATED WORK

notified when the job execution finishes. The processes are defined at the OpenEO core API and therefore independent from the back end they get called, other than the data identifier, which is different for every back end.

The previous example shows a process graph that only uses predefined processes and data. Within the OpenEO project there is the possibility to define individual processes and execute them on the back end. In the project they are called “user defined functions” and are at the writing of this thesis not well defined, but are basically code written by the OpenEO user that gets sent to the back end and executed at a secure environment. The user can define processes and can run them with the data provided at the back end, using the infrastructure of the back end. Every back end has to individually define what the restrictions on user defined functions are. There is also the possibility to upload files to the back end that can be used within the process graph. The user can upload the file through an explicit endpoint of the back end. After the upload finished, the file can be referenced inside the process graph.

2.4.2 Back end Overview

Even though the back ends implement the OpenEO core API standard, they are still very diverse behind the abstraction layer. Some back ends has already an API, where the OpenEO calls have to be adapted to. There are 7 partners within the OpenEO project that are implementing a back end driver. The back ends have to manage the translation of the process graph to the actual code that executes the defined process chain. Also the billing of the users is completely different on every back end.

Add section about
EO reproducibility.

```

{
  "process_id": "min_time",
  "args": {
    "imagery": {
      "process_id": "/user/custom_ndvi",
      "args": {
        "imagery": {
          "process_id": "filter_daterange",
          "args": {
            "imagery": {
              "process_id": "filter_bbox",
              "args": {
                "imagery": {
                  "product_id": "S2_L2A_T32TPS_20M"
                },
                "left": 652000,
                "right": 672000,
                "top": 5161000,
                "bottom": 5181000,
                "srs": "EPSG:32632"
              }
            },
            "from": "2017-01-01",
            "to": "2017-01-31"
          }
        },
        "red": "B04",
        "nir": "B8A"
      }
    }
  }
}

```

Figure 2.2: Process graph example (from Github Spec)

CHAPTER 3

Design

In this chapter the design of the provenance capturing gets described. First the architecture of the OpenEO project gets described further. Especially the job execution path throughout the parts of OpenEO gets described, because it is the key element for the context capturing. Therefore the back end, where the jobs get executed get described in more detail. After that the context model gets presented.

In this section the concept of the provenance capturing is described. The aim of the capturing is to be easy to implement on the OpenEO back ends, but also powerful enough to make the reproduction of the workflows possible. The whole capturing process can be structured in three thematic parts.

1. Back end provenance

This part describes the provenance of the back end, that is not depending on a job execution. There has to be an automated process running on the process machine to capture every change that influences the execution of OpenEO jobs. A description of more detail can be viewed on chapter

2. Job dependent provenance

This part contains the provenance of a job execution. It consists of the workflow context and the data related to one specific job. This in combination of the back end provenance is the full provenance data to make an OpenEO job theoretically re-executable.

3. User information

This part consists of all possibilities that provides users with information they want to know about the provenance of the back end and the provenance of the job execution. The aim is to provide users with information that makes them easier to choose between back ends. On the other hand the user shall be able to retrieve the

provenance data of the job in a way that the back end has no security risks. This part is not directly part of the context model, but of the implementation on the user interface.

3.1 Back end provenance

The scope of this part of the context model is to get the static environment of where the job execution at the back end takes place. It contains the provenance data that is independent from a job execution, so does not change regardless of how much jobs were processed. It can only change from inside of the back ends by their maintainers. In the OpenEO project there are a great variety of different back end providers with very different setups, so the challenge of this part is to make it as simple and generic as possible. The data captured in this part of the context model is not meant to be shown to the user directly, because of security issues. The following data gets suggested to be the minimal set of static provenance data that has to be captured.

1. Github Repository

Since the project is an open source project and every back end has a github repository were at least the basic setup is stored. The aim of this strategy is to get other back ends with similar settings to reuse the already working setups of the running back ends. Therefore, at least during the project runtime, every back end provider has a github repository where the code is publically available. This information is added to the back end provenance by saving the git repository url, the used branch, the used commit and local changes to the repository.

2. Used Folders

The back end will most probably not only use the github repository for the processing, so that other folders are involved that stores config files or additional code. This information is also crucial for the job execution and therefore needs to be stored. Changes of this directory have to be detected and stored in the context model. Temporary written folders have to be excluded to the capturing to prevent detecting changes that are not important.

3. OS and packages

The operating system can also have an impact on the processing and so it has to be stored to the context model too. Not only the information about the operating system, but also the installed packages have to be stored in the context model. In a linux based system it is sufficient to store the location where the packages get installed (e.g. /etc). If the whole processing is done in a virtual container the operating system of that container shall be stored. If the whole processing is done in a docker container, the docker description file has to be saved in addition. But only if the docker container description does not change on different input process graphs, otherwise it has to be added to the “job dependent provenance”.

4. Core API Version

The back ends core api version is the version of the OpenEO core API version it is using. In the productive usage of the OpenEO project it can happen, that the versions of the clients and the back ends differ, so that the behaviour might not be compatible. So the back end server configuration is depending on the core API version of OpenEO, hence it shall be added to the context model.

5. Back End Version

The captured data described in the previous sections will result in a back end version. The back end version shall be an identifier (e.g. a number) that gets updated on every change of the backend considering the provenance data described above.

To provide long term stability of the capturing, there shall be a tool to automatically capture the provenance data to the context model. Every change on any of the previously described context data has to be detected automatically and have to result in a new back end version. This can be done by a standalone capturing tool that does not need to have detailed insights into the back end.

3.2 Job dependent provenance

In this section the job dependent provenance of the context model gets described. The data captured is tied to an specific job execution, so for every job execution a new context model gets created. The structure of the capturing is using the defined process graph structure of OpenEO. The process graph is already a description of the processes that run at the back end. So sending the same process graph to the same back end again shall result in the same outcome. To assure that the process graph results in the same way of processing the processing itself gets captured. After receiving the job the back end will transform the process graph into code that actually does the processing. Figure

The single parts of the capturing are described more in detail in the following sections.

3.2.1 Input Data

The input data of the processing is crucial for the outcome of the process graph. Even though the process graph already uses an identifier for the input data that has to be unique for the specific back end, changes to this data might not result in a new identifier. So jobs called later in time might use another version of the input data than the previous jobs. To prevent this the input data has to be persisted according to the 14 steps of data provenance defined by the RDA[1] described in the related work section. Every back end is responsible for the data versioning, for this context model it is assumed that the back ends preserve the data properly. Nevertheless the current core API version (v0.0.2) is not capable of letting the user choose, which version of data the user want to use. Hence there has to be an addition to the process graph so that product identifier shall have an additional timestamp. The meaning of the timestamp is to tell the back end that

the user wants to use the data at the time of the timestamp. This gives the back ends the possibility to design their own data versioning. They just have to remember, which version of the data was available on what date and time. Not only the product used for the process graph is input data but also files uploaded by the user are included at the input data. As long as the file name is accessed at the specific process graph the file gets added to the provenance of the job execution.

To guarantee that the data has not changed even though the version has not changed it is considered that a hash of the input data is added to the provenance data. The earth observation data stored at the OpenEO back ends can be very big and therefore it is recommended that the back end provider already hash the input data after every change on it so that the performance of the job execution is not threatened. This also applies to the uploaded files of the user, where the hash can be done already after the uploading. In OpenEO files that get uploaded with the same filename are automatically overwritten so that it is especially important to take the hash of this files for the context model. Otherwise it is not possible to replicate differences on a job execution.

3.2.2 Process Data

In this section the capturing of the process itself gets described. As mentioned above the process is described at the process graph, but to be certain that the same process id results in the same code, the code has to be captured. So the code running the processing has to be captured and has to result in a signature that can be compared to other executions. An example of doing so is to hash the entire resulting code. The way of executing a specific process graph is not only related to the code running it, but also by the dependencies of the code. So the environment of the code and information about the code environment has to be added to the context model. If there are any docker container, where the job is running that is not independent on the process graph has to be captured. The programming language, its version and its used packages has also to be captured and added to the job context model.

3.2.3 Output Data

The output data of the processing has to be captured as well. In the OpenEO project the results can be rather big, that is why a simple generated identifier of the output result is enough in this context model. An example of this identifier would be a hash over the resulting data. The aim of the output data capturing is not to add the possibility to find differences between results, but to be able to see that the results are different. So for this context model a simple solution is sufficient.

3.2.4 Benchmarking

So far only the capturing of a single process, including input and output data, are described, but not how the whole process graph shall be captured. The basic idea of the job dependent provenance capturing is to capture the input data of the process

graph, the whole process graph and the output of the process graph, described like in the sections above. This is a common basis the back end providers can agree on and is not affecting the back end providers implementations much. It is also capable of giving the OpenEO user a simple overview of how the results are generated and what are differences between different job executions. To make the comparison of different execution behaviors more informative, the granularity of the capturing can be improved. So that is how the benchmarking context model gets introduced. The idea is that not only the input data the whole process and the output data gets captured, but also every data in between of every process in the process graph. So for every process in the process graph the input, output data gets captured and the code executing the process. This makes it easier for users to see were the execution actually changed. But it also comes with higher implementation costs at the back end side. It also will affect the performance of the execution more than the common context model. The possibility of the implementation of such a granularity is also highly dependent on the back end implementation. Not every back end might be able to implement this into the back end due to execution optimization where the execution and results of the single processes in the process graph are not distinctable. Nevertheless in theory the granularity does not have to stop there, the data and code could be captured for every line in the code. This flexibility of capturing makes it also a context model that can be adapted to different use cases in the future.

3.3 User Information

Read through
above text and
improve.

Add section about
noworkflow



Implementation

- 4.1 The EODC Backend
- 4.2 Job Independent Context
- 4.3 Job Dependent Context
- 4.4 User Interface
- 4.5 Overview of Changes

Enter your text here.

CHAPTER 5

Evaluation

5.1 Use Cases

5.2 Using Python Client

CHAPTER 6

Conclusion and future Work

Enter your text here.

List of Figures

2.1	Overview of the OpenEO architecture. (from openeo.org)	5
2.2	Process graph example (from Github Spec)	7

List of Tables

Index

distribution, 5

Glossary

editor A text editor is a type of program used for editing plain text files.. 5

Acronyms

CTAN Comprehensive TeX Archive Network. 11

FAQ Frequently Asked Questions. 11

PDF Portable Document Format. 6, 10, 11, 15

SVN Subversion. 10

WYSIWYG What You See Is What You Get. 9

Bibliography

- [MGC⁺15] Luc Moreau, Paul Groth, James Cheney, Timothy Lebo, and Simon Miles. The rationale of prov. *Web Semantics: Science, Services and Agents on the World Wide Web*, 35:235 – 257, 2015.
- [MPM⁺13] Tomasz Miksa, Stefan Pröll, Rudolf Mayer, Stephan Strodl, Ricardo Vieira, José Barateiro, and Andreas Rauber. Framework for verification of preserved and redeployed processes. In *Proceedings of the 10th International Conference on Digital Preservation, iPRES 2013, Lisbon, Portugal, September 2 - 6, 2013*, 2013.
- [PWS⁺17] Edzer Pebesma, Wolfgang Wagner, Matthias Schramm, Alexandra Von Beringe, Christoph Paulik, Markus Neteler, Johannes Reiche, Jan Verbesselt, Jeroen Dries, Erwin Goor, and et al. Openeo - a common, open source interface between earth observation data infrastructures and front-end applications. Nov 2017.
- [RAvUP16] Andreas Rauber, Ari Asmi, Dieter van Uytvanck, and Stefan Pröll. Identification of reproducible subsets for data citation, sharing and re-use. 2016.
- [RMD12] H. Ramapriyan, J. Moses, and R. Duerr. Preservation of data for earth system science - towards a content standard. In *2012 IEEE International Geoscience and Remote Sensing Symposium*, pages 5304–5307, July 2012.
- [RMH⁺11] David De Roure, José Manuel, Kristina Hettne, Khalid Belhajjame, Raúl Palma, Graham Klyne, Paolo Missier, José Enrique Ruiz, and Carole Goble. Towards the preservation of scientific workflows. In *In Procs. of the 8th International Conference on Preservation of Digital Objects (iPRES 2011)*. ACM, 2011.
- [RMMP15] Andreas Rauber, Tomasz Miksa, Rudolf Mayer, and Stefan Pröll. Repeatability and re-usability in scientific processes: Process context, data identification and verification. In *DAMDID/RCDL*, 2015.