

The gocfl Workflow for Research Data Reporting and publication of clustered data packages

The Oxford Common File Layout (OCFL¹) is an application-independent specification for the structured, transparent storage of digital information. Its internal structure is geared towards the requirements of long-term storage and makes it easier to check the completeness and integrity of OCFL-data (-packages). The main goal is to keep knowledge about data available, even after active use, e.g. after completion of a research project.

The open-source software GOCFL² is programmed in GO and fulfills the OCFL standard. It provides a routine for the automated and sustainable packaging of data. The tool enables resource-saving reporting. In addition, people who were not actively involved in the editing processes can also keep an eye on the data situation (content, usage and other property rights) with the help of the web- and pdf-based human-readable reports. Optionally, both the OCFL packages and the reports (web format + PDF) can be published FAIRly.

This document outlines the implementation of the GOCFL workflows at the Academy of Art and Design Basel FHNW. It explains the core-functionality of this specific open-source software. Graphically, the GOCFL workflow can be anchored in the research data cycle as follows:

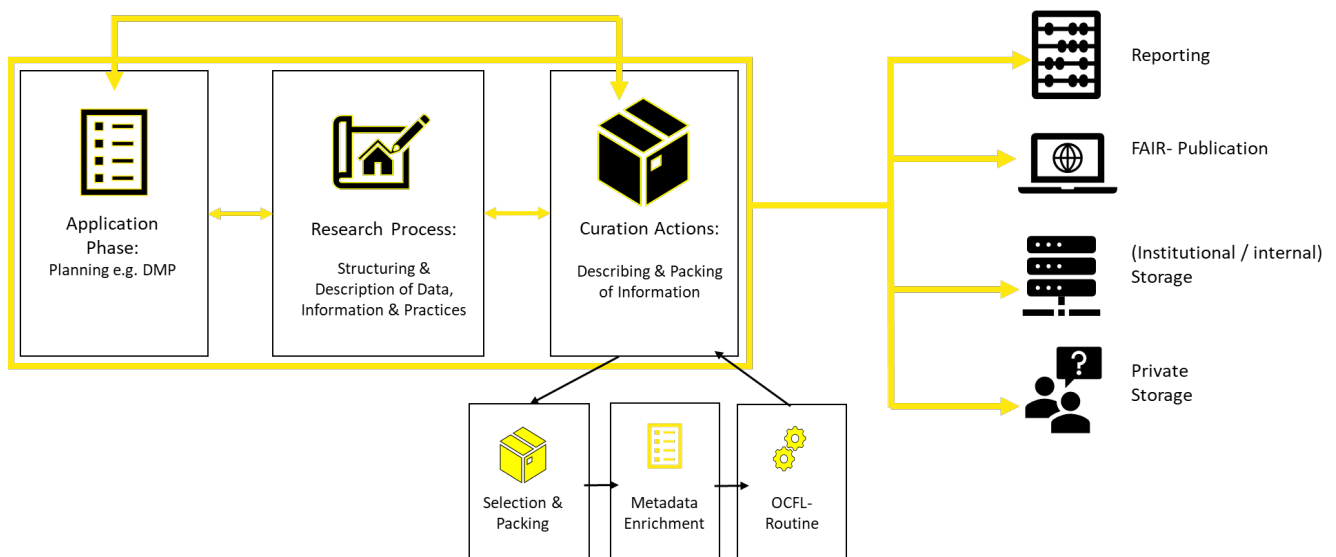


Fig. 1: GOCFL workflow as implemented within the research data cycles at HGK Basel FHNW

¹ <https://ocfl.io/1.1/spec/>

² <https://github.com/je4/gocfl>

Implementing of GOCFL at HGK Basel FHNW

The open-source software GOCFL was integrated within the research data cycle in 2023 at HGK Basel FHNW. The main aim was to improve both the quality of data management regarding data overview and the publication options for closed data. Since then, GOCFL has been one of the institutionally funded services used in the final research phase when a project is completed (see Fig. 1).

A classic application for the use of GOCFL is research data that should not or cannot be freely published after completion of a project. The reports comment on data protection, copyright, and personal rights. They offer researchers an efficient way of representing the diversity of material identified in the data management plans, while they enable research managers to obtain an overview of the status quo of the data available at their institute.

Functionalities of GOCFL

As a software library, GOCFL extends the OCFL standard according to specified rules with software routines that contribute to sustainability, reporting and thus improve FAIR-ness. GOCFL zips folders of structured data and preserves the files and their structure. It extracts technical metadata, builds checksums and documents the applied routines. GOCFL allows later updates and validation through versioning. Of special interest are the reports GOCFL creates: They provide clear and human-readable information about the available data packages as browser-based (full length) and PDF-compatible (short/abbreviated) versions.

Implemented as a modular software with indexing cascades, further functionalities can be added or removed to meet specific requirements. To give an example, if required, GOCFL might encrypt filenames.

In general, all OCFL objects contain a so-called object root, in which data and an inventory are contained.

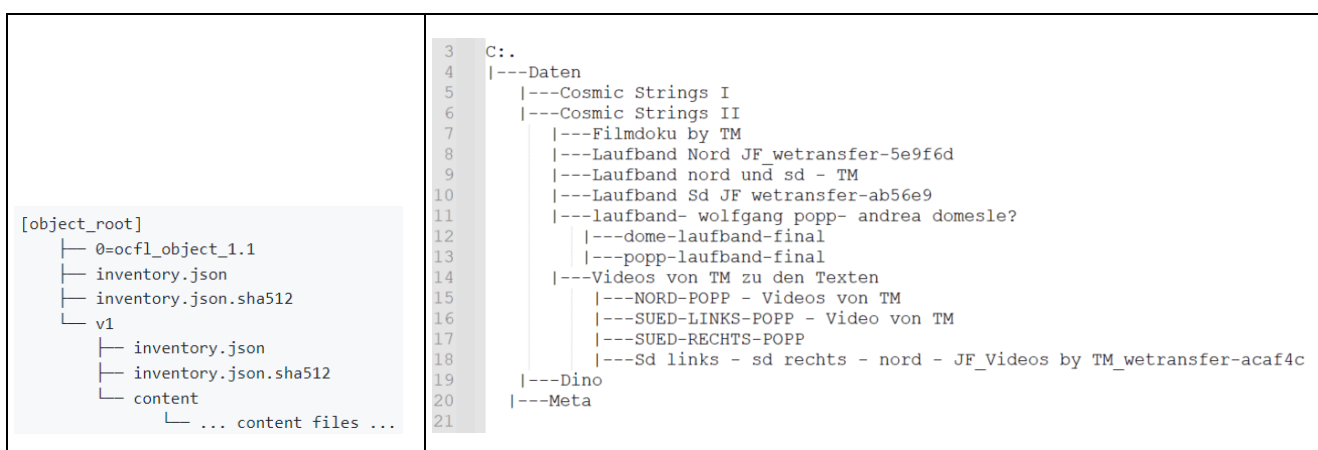


Fig. 2: OCFL object structure: Left side according to the standard³, right as sample of HGK Basel

Since filenames and filing structures often contain semantic information regarding content and context such as the research process (dates, persons or places involved), filenames and -structures are preserved. Corresponding data is extracted and can be displayed after the routine.

³ <https://ocfl.io/1.1/spec/>

Since special characters in file names can jeopardize sustainability, one of the GOCFL-routines re-names them by converting unsupported characters to utf8. All changes are documented.

Due to its technical structure as cascade of single software libraries, which can be added or removed as wished, GOCFL can be extended for carrying out additional requirements.

Currently GOCFL extracts the following metadata during the technical analysis:

- File size
- Creation and modification date
- Filename and position in the directory structure
- Characterization of the data type and the underlying operating system
- Determination of the file format using PRONOM Persistent Unique Identifier (PUID)
- Creation of checksums to identify the files
- Depending on the configuration, further metadata
- Removal of special characters from file names to increase sustainability
- Creation of a ZIP file⁴.

To create an easy-to-understand report, GOCFL also offers the following functionalities:

- Statistical evaluation of the data according to data type and format
- Creation of thumbnails
- Creation of a web-based overall report
- Creation of a short version of the report that is laid out for common PDF format
- Filename and position in the directory structure
- Characterization of the data type and the underlying operating system

⁴ The OCFL packages are available as a zip file, as archive memories are generally not very good at handling many small files.

Web based report

The index and technical metadata of the OCFL object can be viewed using the browser:

gocfl FHNW-HGK-MED_Videocity_Eller

Object
Manifest
Version v1
Version v2
Report

Versions

v1: Tabea Lurk	v2: gocfl
2023-09-06 12:51:02 +0200 CEST mailto:tabea.lurk@fhnw.ch initial add	2023-09-06 12:55:50 +0200 CEST https://github.com/jk4/gocfl automated version

Statistics

Quantities

Number of Files: 109
Number of different Files: 91
Size of Data: 11 GB
Files without size data: 6

Pronoms

- 2 fmt/111
- 1 fmt/1149
- 40 fmt/199
- 1 fmt/276
- 2 fmt/353
- 1 fmt/394
- 14 fmt/41
- 19 fmt/503
- 1 fmt/645
- 1 x-fmt/384
- 1 x-fmt/391

Mimetypes

- 40 application/mp4
- 1 application/octet-stream
- 1 application/pdf
- 2 application/x-tika-msoffice
- 16 image/jpeg
- 3 image/tiff
- 12 multipart/appledouble
- 1 text/markdown
- 1 video/quicktime

Fig. 3: Web-based report of the packed OCFL object

Manifest FHNW-HGK-MED_Videocity_Eller

Object
Manifest
Version v1
Version v2
Report

- v1/content/README.md (111 B)
- v1/content/data/Cosmic=u0020Strings=u0020/c66_Tomas=u0020Eller.tiff (1.3 MB)
- v1/content/data/Cosmic=u0020Strings=u0020/cosmicstrings-eller-07.jpg (442 KB)
- v1/content/data/Cosmic=u0020Strings=u0020/cosmicstrings-eller-08.jpg (560 KB)
- v1/content/data/Cosmic=u0020Strings=u0020/cosmicstrings-eller-09.jpg (540 KB)
- v1/content/data/Cosmic=u0020Strings=u0020/cosmicstrings-tomaseller-basel_1.mp4 (1.6 GB)
- v1/content/data/Cosmic=u0020Strings=u0020/|Filmdoku=u0020by=u0020Tomas=u0020Eller/|DS_Store (6.1 KB)
- v1/content/data/Cosmic=u0020Strings=u0020/|Filmdoku=u0020by=u0020Tomas=u0020Eller/|_DS_Store (6.1 KB)
- v1/content/data/Cosmic=u0020Strings=u0020/|Filmdoku=u0020by=u0020Tomas=u0020Eller/|COSMICSTRINGS-BASEL-2018.mp4 (1.6 MB)
- v1/content/data/Cosmic=u0020Strings=u0020/|Filmdoku=u0020by=u0020Tomas=u0020Eller/|COSMICSTRINGS-BASEL-2018.mp4 (1.6 MB)
- v1/content/data/Cosmic=u0020Strings=u0020/|Aufband=u0020Nord=u0020Julia=u0020Franck_wetransfer-5e9fd/|_01_NORD-JULIAFRANCK.mp4 (60 MB)
- v1/content/data/Cosmic=u0020Strings=u0020/|Aufband=u0020Nord=u0020Julia=u0020Franck_wetransfer-5e9fd/|01_NORD-JULIAFRANCK.mp4 (144 MB)

Fig. 4: Manifest of the package with list of files, linkage to the full index information, file format and size information

Checksums

sha512	ae59db15a8decbaf2843d7cd45d56446659c3fb4a1531c0227210bc9a89af33ed088952ad05496f1deecb3d9eeebf5478f26b5dcae49a1a7123a8984446a
blake2b-384	a587d7f3395b6b594879d20e5278194d3b1e156df9e1b2bc386c501c3394454b6c4154b7802881ad53fa2b3819
md5	453d7e007fb1ca2e70534eb6a7a5d96c
sha256	393ac42a1d4e916eaa16d3cb47600cf005de969a909bb5ec5d716c335054313

Internal Names

v1/content/data/Cosmic+u0020strings+u0020/cosmicstrings-eller-07.jpg

External Names

Version	Name	Time	Size	Attr	OS	Sys
v1	data\Cosmic Strings \cosmicstrings-eller-07.jpg	C: 2023-09-06 12:39:35 A: 2023-09-06 12:51:03 M: 2018-05-30 10:09:10	463 kB	Archive	windows	<pre> { "creationTime": { "HighDateTime": 31855846, "LowDateTime": 184885539 }, "fileAttributes": 32, "fileSizeHigh": 472777, "lastAccessTime": { "HighDateTime": 31956848, "LowDateTime": 131686583 }, "lastWriteTime": { "HighDateTime": 30668781, "LowDateTime": 209681382 } } </pre>
v2	data\Cosmic Strings \cosmicstrings-eller-07.jpg					

Thumbnail

Thumbnail Process: image#01

Image




Fig. 5: Excerpt of the detail view of the analysis, including a thumbnail of the stored images

PDF-based report

The PDF report contains a laid-out extract of the technically extracted metadata. The cover page with the info.json (see below) is followed by a statistical list of the file formats contained in the package, the directory tree and finally a short version of the extracted metadata.

Report

mediathek:FHNW-HGK-SNF_185436-b

Der vorliegende Report basiert auf der ODFL-Dokumentation des Datenpakets subFSzipSfjFiles(osFSRWIC:comp/lati)/Exp/Journeys.zip(j)/mediathok_FHNW-HGK-SNF_185436-b). Er enthält einen Auszug der Metadaten, die im Zuge der Archivierung für jede der im Datenpaket enthaltenen Dateien generiert wurden.

This report is based on the ODFL documentation of the data package subFSzipSfjFiles(osFSRWIC:comp/lati)/Exp/Journeys.zip(j)/mediathok_FHNW-HGK-SNF_185436-b). It contains an excerpt of the metadata generated in the course of archiving for each of the files contained in the data package.

OCFL-Object Version v1
 Created: 2023-06-20 16:28:58
 Name: Lark, Tubas
 Address: mailto:tubas.lark@fhnw.ch
 Message: Initial commit

generated by goctl / https://github.com/jel/goctl

MEDIATHEK:FHNW-HGK-SNF_185436-B

Technology – Human – Design: Paradigms of Ubiquitous Computing - Journey

Hochschule für Gestaltung und Kunst FHNW Basel
IXDM Research Projects

Torpus, Jan
mailto:info@ixdm.hgk@fhnw.ch
2023-06-20

The research setting *Journey* focuses on human behavioral patterns and appropriation processes in unknown techno-social hybrid environments. We investigate and compare both perspectives: user experience and machine recognition of human behavior. To do so, we composed an interactive research setting of six separate spaces to stage paradigms of ubiComp such as: ubiquity, immediacy, invisibility, seamlessness/seamfulness and interconnectedness. By applying strategies and aesthetics from the media arts and by implementing technological sensor-actor systems, we created a measurable and debatable environment. It offers snippets of narratives that can be mentally combined and is carefully designed to enable test persons to fully engage and willingly suspend their disbelief. The bodily interference with the physical environment is intensified by an interactive garment that test persons put on before their journey. By doing so, they leave their familiar comfort zone and take on a new role in the narrative.

Alternative Title	UbiCombs Journey
Signature	FHNW-HGK-SNF_185436-b
Organisation	FHNW-HGK
ID	IXDM-Research
Sets	IXDM_HGK_FHNW_SNF_Forschung
Keywords	Arts, Mediaepistemologie, Sensor-Actor-Mapping, Internet der Dinge, Responsive Environment, Sensor-Network, Ubiquitous Computing, Kunstforschung, Data Mapping
Identifiers	oai:2.1641719.S9EUAJIS https://ubicombs.ch/journeys/
References	FHNW-HGK-SNF_185436-a
Deprecates	
Ingest Workflow	WF01

Page 2 of 145

generated by goctl / https://github.com/jel/goctl

Fig. 6: Intro sequence of the PDF-based report (page 1) with information about the project and persons in charge (page 2)

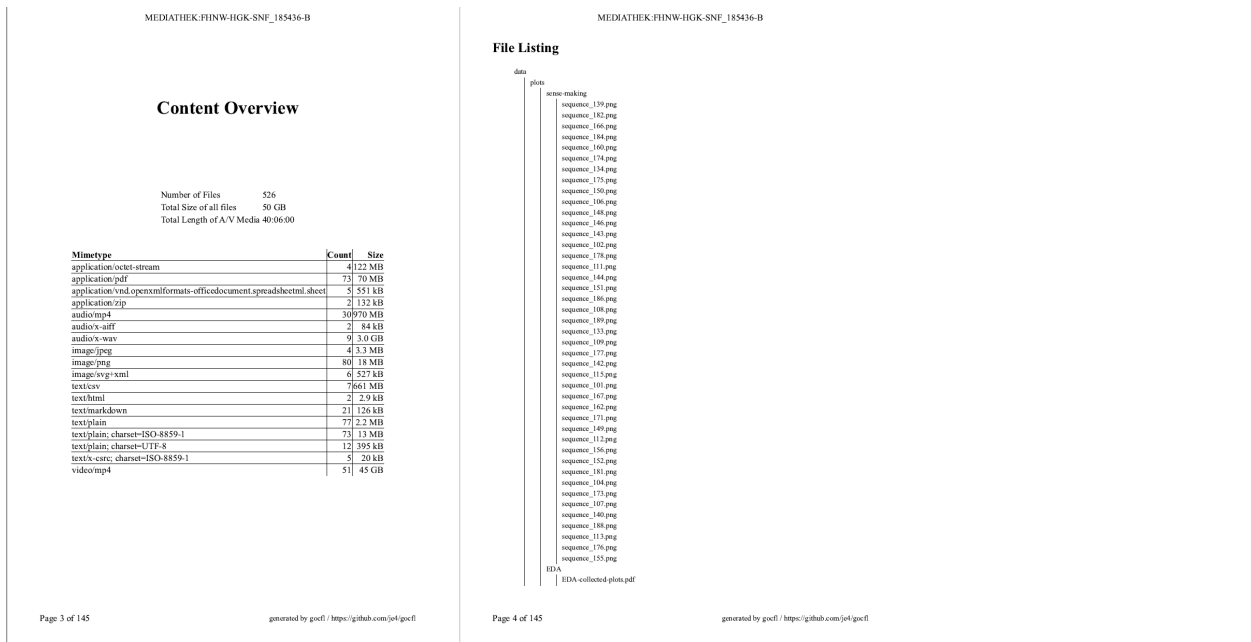


Fig. 7: Listing of file formats according to frequency and file tree

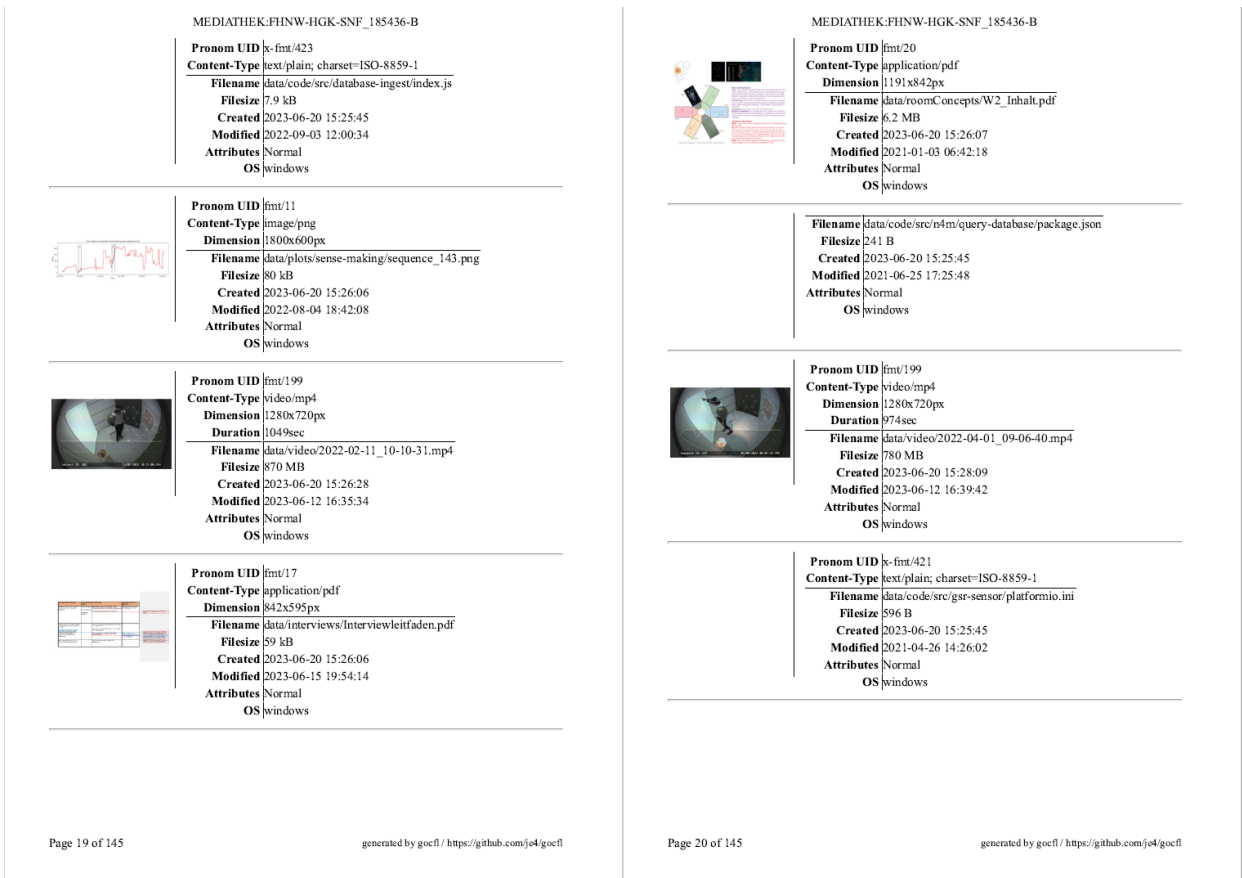


Fig. 8: Sequence with thumbnails and metadata information about the file such as PUID

Workflow

The GOCFL workflow requires the following steps:

- Curation of data
 - Store files in one or more folders
 - Describe/characterize data clusters with a readme file (see below)
 - Add basic record information in an info.json file (see below)
- Configure GOCFL (once per executing instance [computer/server])
- Start help tools (Tika, Powershell)
- Execute the GOCFL script
- (Re-)View report and archive package and forward to a publication instance if necessary.

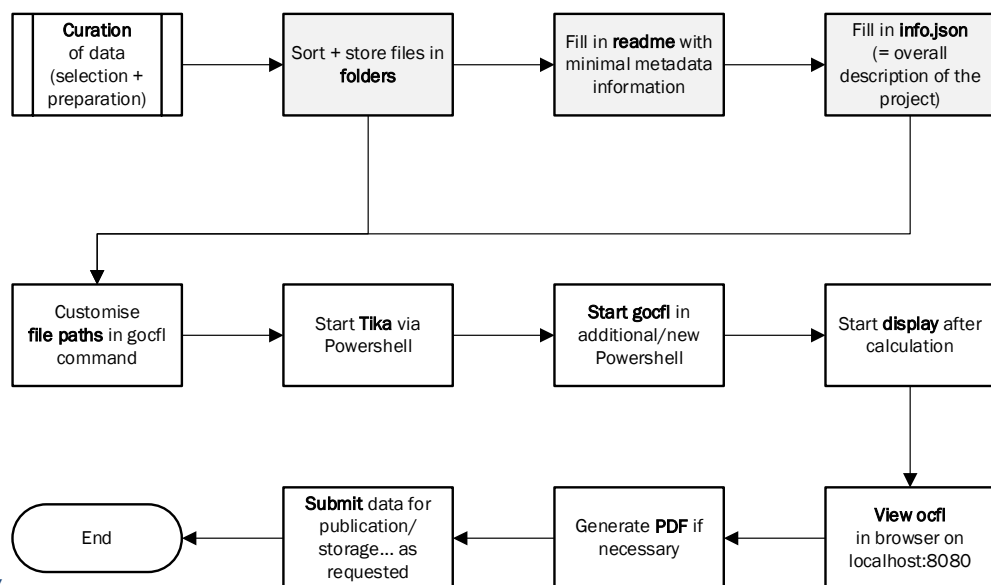


Fig. 9: OCFL-workflow

Description

Researchers are responsible for determining the depth, quality and format of description. The following explanation may thus be regarded as recommendations for a minimum amount of information, still increasing the FAIRness of data.

Nevertheless, we recommend creating at least one readme file per superordinate folder/file path that provides information about the contents of the directory. If different legal areas are affected, these should be mentioned. In principle, two levels of description should be used:

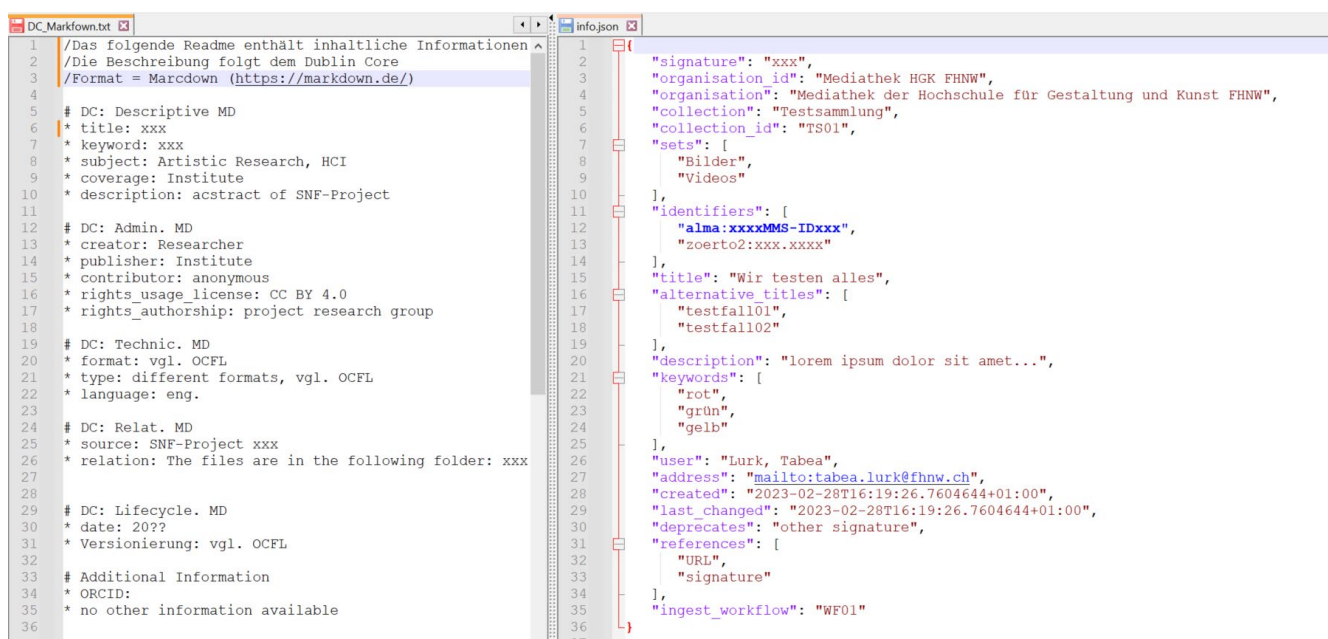
- Superordinate information about the project, the persons responsible and the creation --> info.json
- Content characterization of the stored files --> readme.

info.json

The requirements for info.json (Fig. 10) are stricter, as they must be generally meaningful. They are prominently placed on the cover page of the report and should contain formal information such as the title of the project, responsibilities and legal issues beyond the active research period.

Readme

For the readme (Fig. 10), we recommend structuring the data in terms of content according to Dublin Core. Formatting conventions vary by nature. A suitable storage format might be raw text, containing e. g. a Marcdown structure, or JSON.



```
DC_Markdown.txt infojson
1 /Das folgende Readme enthält inhaltliche Informationen
2 /Die Beschreibung folgt dem Dublin Core
3 /Format = Markdown (https://markdown.de/)
4
5 # DC: Descriptive MD
6 * title: xxx
7 * keyword: xxx
8 * subject: Artistic Research, HCI
9 * coverage: Institute
10 * description: acstract of SNF-Project
11
12 # DC: Admin. MD
13 * creator: Researcher
14 * publisher: Institute
15 * contributor: anonymous
16 * rights_usage_license: CC BY 4.0
17 * rights_authorship: project research group
18
19 # DC: Technic. MD
20 * format: vgl. OCFL
21 * type: different formats, vgl. OCFL
22 * language: eng.
23
24 # DC: Relat. MD
25 * source: SNF-Project xxx
26 * relation: The files are in the following folder: xxx
27
28
29 # DC: Lifecycle. MD
30 * date: 20??
31 * Versionierung: vgl. OCFL
32
33 # Additional Information
34 * ORCID:
35 * no other information available
36
37
38 {
39   "signature": "xxx",
40   "organisation id": "Mediathek HGK FHNW",
41   "organisation": "Mediathek der Hochschule für Gestaltung und Kunst FHNW",
42   "collection": "Testsammlung",
43   "collection_id": "TS01",
44   "sets": [
45     "Bilder",
46     "Videos"
47   ],
48   "identifiers": [
49     "alma:xxxxMMS-IDxxx",
50     "zoerto2:xxx.xxxx"
51   ],
52   "title": "Wir testen alles",
53   "alternative_titles": [
54     "testfall01",
55     "testfall02"
56   ],
57   "description": "lorem ipsum dolor sit amet...",
58   "keywords": [
59     "rot",
60     "grün",
61     "gelb"
62   ],
63   "user": "Lurk, Tabea",
64   "address": "mailto:tabea.lurk@fhnw.ch",
65   "created": "2023-02-28T16:19:26.7604644+01:00",
66   "last_changed": "2023-02-28T16:19:26.7604644+01:00",
67   "deprecates": "other signature",
68   "references": [
69     "URL",
70     "signature"
71   ],
72   "ingest_workflow": "WF01"
73 }
```

Fig. 10: Suggestions: left - content description of data, right - formal information of context

At the media library of the HGK Basel FHNW, GOCFL is mainly used in the areas of research support and for clarifying copyrights in the course of the data transfer for special collections.

License

GOCFL was developed by Dipl. Inform. Jürgen Enge⁵ and is released with Apache-2.0 License.

⁵ <https://ub.unibas.ch/de/kontakt/>.