

## DeepGreen - ein schnittstellengetriebener Ansatz zur Datenanreicherung und -analyse: Konzeption eines Workflows mit Modulen in RStudio



111 BiblioCon2023 Hannover 24.05.2023 DeepGreen Anwender:innen-Treffen – Austausch und Weiterentwicklung

煎茶 (Sencha) - 八女市 (Yame-shi) 福岡県  
(Fukuoka-ken) Juli 2015 T.Wetzenstein

## Inhalte

- Charakteristika der DeepGreen-Sammlung auf MACAU
- Datenschnittstellen
- Anreicherung der DeepGreen-Daten in MACAU
- Wege der Datenanalyse: Harvesten, Transformation, Filtern, Anreichern, Analysieren
- Sichtbarkeit der DeepGreen-Sammlung auf MACAU
- Nutzen einer verteilten Repositorylandschaft mit neuen Werkzeugen und Schnittstellen für DeepGreen-Sammlungen

## Charakteristika der Sammlung auf MACAU - Workflows

Stand 27.04.2023: **8.080** Dokumente gesamt

5520 Dissertationen

1905 Zweitveröffentlichungen, davon **1011** via DeepGreen

354 Erstveröffentlichungen

304 Verlagsveröffentlichungen (inklusive Kapitel und Artikel)

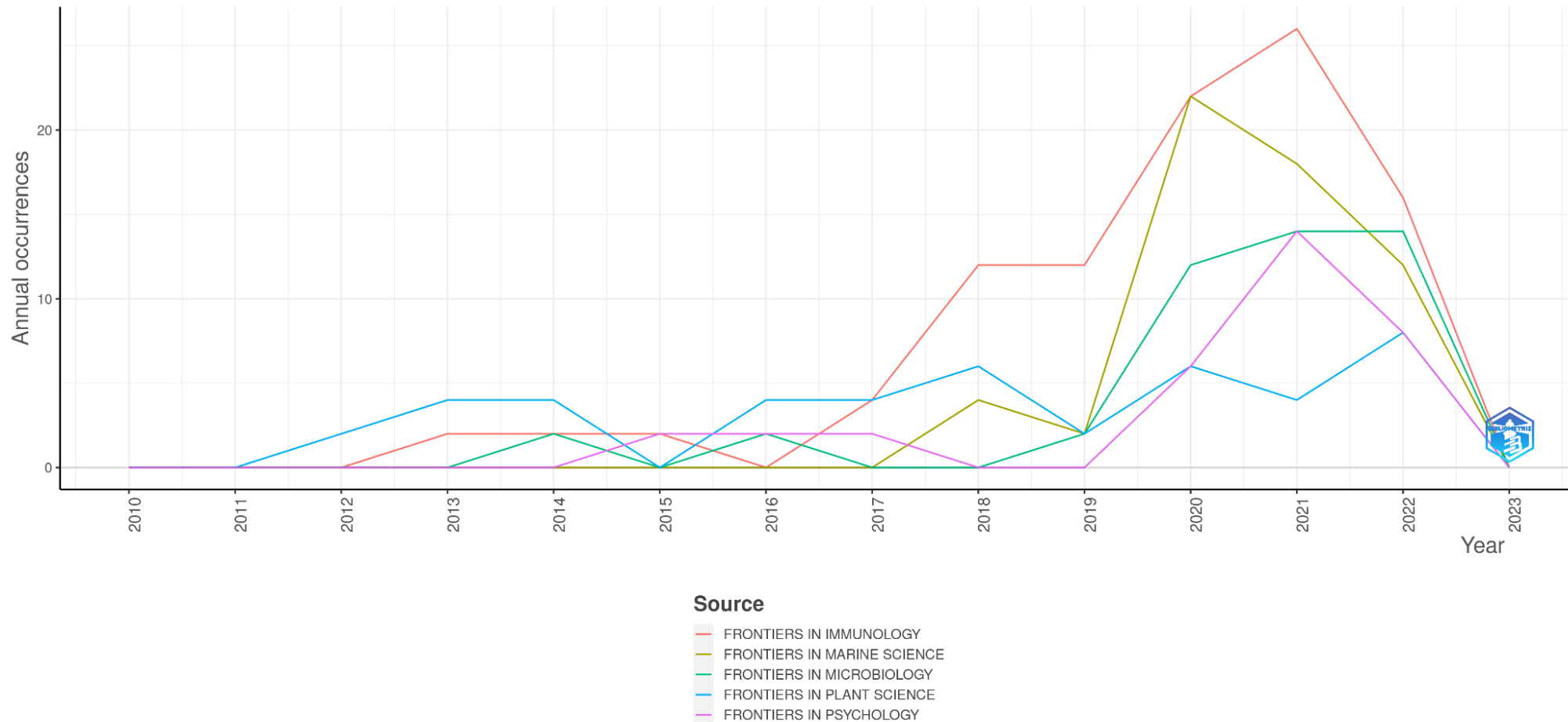
*davon 21 eigenständige Veröffentlichungen (Sammelbände, Zeitschriftenhefte)*

### DeepGreen

|                 |                 |
|-----------------|-----------------|
| <b>872</b> Dok. | CC-BY 4.0       |
| 35 Dok.         | CC-BY-NC 4.0    |
| 23 Dok.         | CC-BY-NC-ND 4.0 |
| 20 Dok.         | CC-BY 3.0       |
| 3 Dok.          | CC-BY 2.0       |
| <b>53</b> Dok.  | Ohne CC-Lizenz  |

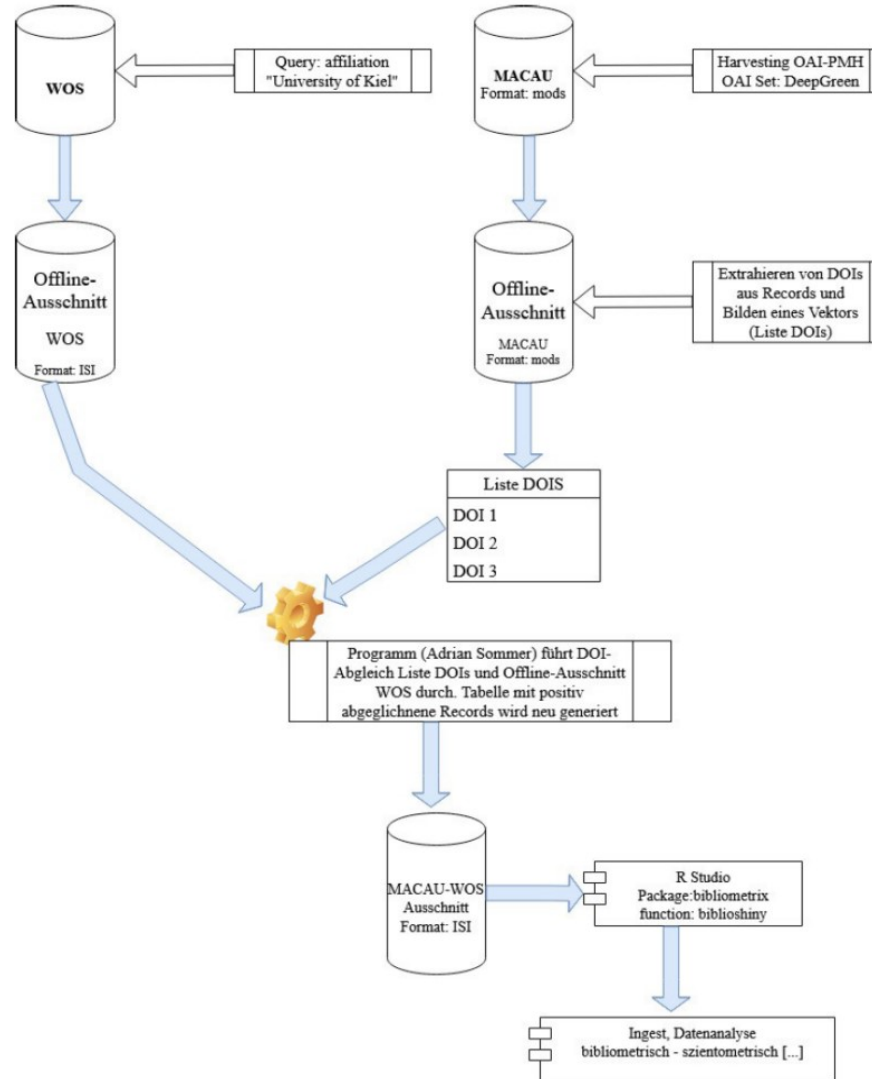
## Charakteristika der DeepGreen-Sammlung auf MACAU - Produktivität nach Zeitschriften

## Sources' Production over Time



# Wege der Datenanalyse: Analysieren - Inhalte

Durchführung des Workflows  
*realisiert*



## Datenschnittstellen – Web of Science Daten (WOS) – Open Alex – Crossref – DeepGreen

|                                  | JATS ( <u>Publ. Router DeepGreen</u> ) | Crossref (Verlage / <u>Repo.-Betr. / Bibliotheken</u> ) | <u>OpenAlex</u> | WOS ( <u>Clarivate</u> )   |
|----------------------------------|--|---|-----------------|----------------------------|
| <u>Funding Information</u>       | (ja)                                   | ja  | nein            | ja                         |
| <u>Author ORCID</u>              | ja                                     | ja  | ja              | ja                         |
| <u>Author position</u>           | (ja)                                   | ja  | ja              | ja                         |
| <u>ror ID</u>                    | (ja)                                   | ja  | ja              | nein                       |
| weitere Hosts                    | (nein)                                 | (nein)  | ja              | nein                       |
| <u>zitierte Ref.</u>             | (nein)                                 | ja  | ja              | ja                         |
| <u>zitierende Ref.</u>           | (nein)                                 | ja  | ja              | ja                         |
| Lizenzinfo CC                    | ja                                     | ja  | ja              | nein (OA <u>gold: ja</u> ) |
| <u>Ref. Klassifikationssyst.</u> | (ja)                                   | nein  | ja              | ja                         |

MACAU als MyCoRe IR verwendet Datenformat MODS. Bezug der Metadaten durch Resolving der DOIs und interner Verwendung des sog. MyCoRe Enrichment Resolvers

OpenAlex integriert unpaywall-Daten: von dort kommen die Angaben über *alternate\_host\_venue* bzw. *locations*.

Web of Science API war zum gegebenen Zeitpunkt nicht genutzt worden: Daten von Web GUI

## Wege der Datenanalyse: Harvesten - DOI-Liste erstellen

OAI-Harvesting von MACAU vom OAI Set „deepgreen“. Aufruf im Terminal



```
oai-harvest -set deepgreen -metadataPrefix mods https://macau.uni-kiel.de
```

Zählen der Records

```
ls -la | wc -l
```

Anzeigen aller XML Elemente <mods:identifier type="doi">

```
cat *.xml | grep doi
```

Anzeigen aller DOI als Textinhalt des XML-Elements als String

```
cat *.xml | grep doi | cut -d ">" -f2 | cut -d "<" -f1
```

Schreiben der Liste in eine CSV Datei

```
cat *.xml | grep doi | cut -d ">" -f2 | cut -d "<" -f1 > list_doi.csv
```



## Wege der Datenanalyse: Analysieren - Inhalte



- Analysewerkzeuge: CRAN Package bibliometrix
- Data Extraction Tool: CRAN Package openalexR
- Data Extraction Tool: CRAN Package wosr
- Conversion Tool Bibliographische Formate: CRAN Package rbibutils
- Bezug von Daten über Web of Science Webinterface (WOS)





# Wege der Datenanalyse: Analysieren - Inhalte

## Analysewerkzeuge: CRAN Package bibliometrix

127.0.0.1:4836

Most Visited Getting Started Other Bookmarks

**bibliometrix**

biblioshiny

Data

Filters

Overview

- Main Information
- Annual Scientific Production
- Average Citations per Year
- Three-Field Plot

Sources

Authors

Documents

Clustering

Conceptual Structure

Intellectual Structure

Social Structure

|  | TI | U1 | U2 | UT                  | VL  | WC                                   | WE   | Z9 | DB  | AU_UN  |
|--|----|----|----|---------------------|-----|--------------------------------------|--|----|-----|--|
| ASSOCIATION BETWEEN SKIN AND JOINT INVOLVEMENT IN PATIENTS WITH PSORIATIC ARTHRITIS TREATED WITH ADALIMUMAB: ANALYSIS OF DATA FROM A GERMAN NON-INTERV |    | 0  | 2  | WOS:000351253700005 | 230 | DERMATOLOGY                          | SCIENCE CITATION INDEX EXPANDED (SCI-EXPANDED)   | 4  | ISI | GERMANY;GOETHE UNIV FRANKFURT  |
| A QUICK AND UNIVERSAL METHOD FOR STEREOTACTIC VISUALIZATION OF THE SUBTHALAMIC NUCLEUS BEFORE AND AFTER IMPLANTATION OF DEEP BRAIN STIMULATION ELECTRO |    | 1  | 2  | WOS:000188385900016 | 80  | NEUROSCIENCES; NEUROIMAGING; SURGERY | CONFERENCE PROCEEDINGS CITATION INDEX - SCIENCE (CPCI-S); SCIENCE CITATION INDEX EXPANDED (SCI-EXPANDED) | 93 | ISI |  |
| REACHING TREATMENT GOALS IN PSORIASIS WITH CONVENTIONAL SYSTEMIC DRUGS: HOW LONG ARE WE WILLING TO WAIT?   |    | 0  | 1  | WOS:000659539800001 | 238 | DERMATOLOGY                          | SCIENCE CITATION INDEX EXPANDED (SCI-EXPANDED)   | 0  | ISI | UNIV MED CTR HAMBURG EPPENDORF;ABBVIE DEUTSCH GMBH AND CO KG;SANOFI AVENTIS DEUTSCH GMBH;UNIV MED CTR SCHLESWIG HOLSTEIN |

### Options

#### Filters

**Run**

Documents 13885 of 14497  
Sources 1281 of 1310  
Authors 48177 of 48801

#### Language

ENGLISH ENGLISH GERMAN  
ESTONIAN FRENCH GERMAN  
RUSSIAN

#### Publication Year

1966 2023

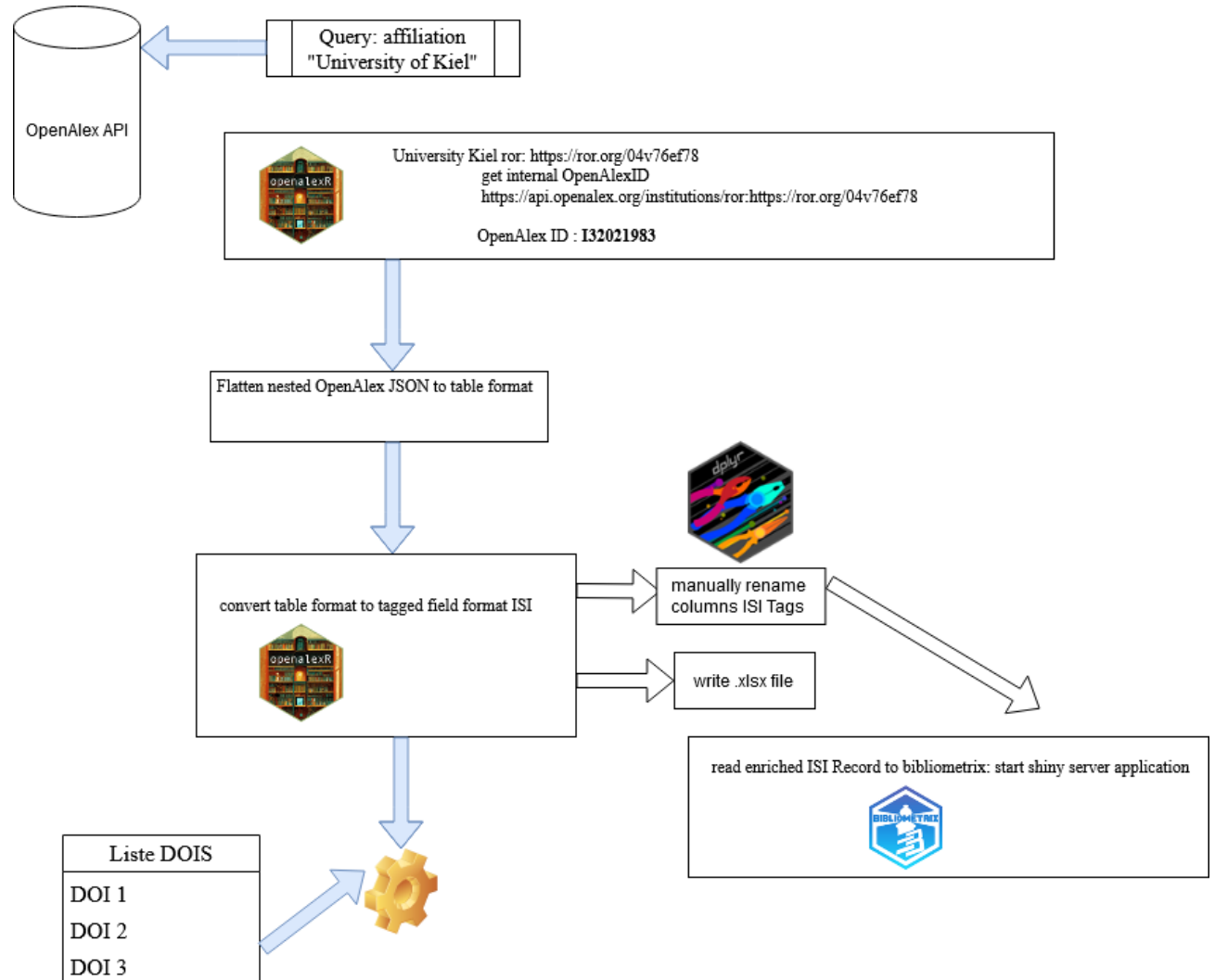
1966 1972 1978 1984 1990 1996 2002 2008 2014 2020 2023

#### Document Type

ABSTRACT OF PUBLISHED ITEM

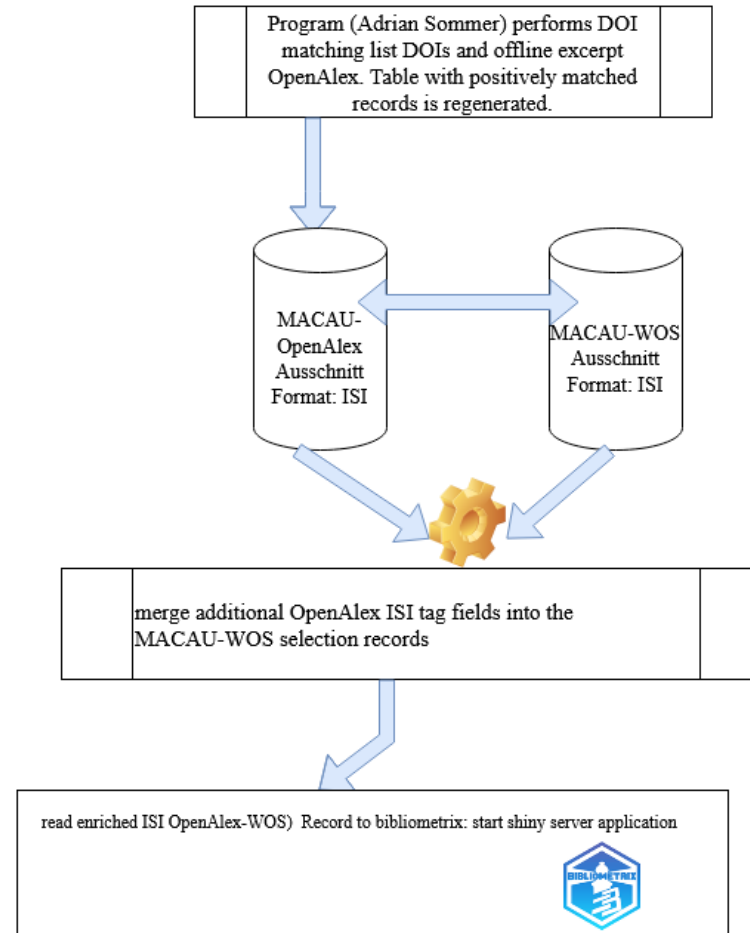
# Wege der Datenanalyse: Analysieren - Inhalte

## Durchführung des Workflows Designvorschlag



## Wege der Datenanalyse: Analysieren - Inhalte

### Durchführung des Workflows *Designvorschlag*



## Anreicherung der DeepGreen-Daten in MACAU

Rstudio Workflow



```
cau_pub_OA <- openalexR::oa_request(query_url = "https://api.openalex.org/works?filter=institutions.id:132021983?", verbose = "TRUE")
```

1. Definieren einer Variablen, in der „works“ der Einrichtung CAU Kiel gespeichert wird. Bezug der Daten von der JSON API OpenAlex (tiefgeschachtelte Daten in JSON Codierung)

```
cau_pub_OA_df <- openalexR::oa2df(cau_pub_OA, entity = "works", verbose = "TRUE")
```

2. Definieren einer Variablen, in der der OpenAlex API Abzug in ein flaches Tabellenformat, sog. Dataframe (df) gebracht wird

```
cau_pub_OA_df_bibmet <- openalexR::oa2bibliometrix(df = cau_pub_OA_df)
```

3. Definieren einer Variablen, in der der Dataframe mit ISI Tags in Tabellenspalten gewandelt wird

In der Funktion `oa2bibliometrix` können Spaltenbezeichnungen teilw. Nicht ISI-gerecht sein, so dass sie in Excel nachbearbeitet werden müssen. Hier liegt evtl. Entwicklungspotenzial für Package-Updates : das erwartete ISI Tag für DOI **DI** ist „doi“ benannt. Ebenso TAG für Schlagwort **DE** als „concept“ – erweiterte Liste siehe Fachaufsatz

## Anreicherung der DeepGreen-Daten in MACAU

Rstudio Workflow



```
WriteXLS::WriteXLS(x = cau_pub_OA_df_bibmet, ExcelFileName = "cau_pub_isi.xlsx", verbose = "TRUE")
```

4. Schreiben des ISI Dataframes in eine Exceldatei

5. Extern oder in Rstudio Umbenennen derjenigen Spalten, für die ISI Tag Bezeichnungen nicht übereinstimmen (Generell sind ISI Tags zum Großteil schon korrekt gesetzt) –  
Steigerung der Transformationsleistung



```
bibliometrix::biblioshiny()
```

6. Einlesen der Exceldatei in Bibliometrix / Auswertungen



## Sichtbarkeit der DeepGreen-Sammlung auf MACAU



MACAU als Datasource auf OpenAlex : ID S4306401923

<https://api.openalex.org/sources/S4306401923>

Kiel University als Institution auf OpenAlex : ID I32021983

<https://explore.openalex.org/institutions/I32021983>

<https://unpaywall.org/sources/repository/ejyjo7aunlatevqewej>

und andere

Serverseitig:

DeepGreen-Sammlung mit eigenem OAI Set `<setSpec>deepgreen</setSpec>`

In Diskussion: HighWire Press Tag `<meta name="citation_collection_id" content="deepgreen">`



OpenDOAR

OpenAIRE | EXPLORE



## Fazit

- Für Mapping in ISI Format von CRAN Package openalexR wären über spätere Updates Verbesserungen zu erwarten für (Erweiterte Liste siehe Fachaufsatz)
- Eine enge Kooperation in der Entwicklung von openalexR und bibliometrix erscheint wichtig vor dem Hintergrund der Formatttransformation nach ISI und die (z.B. in ggplot) vorgesehenen Darstellungen im Package bibliometrix
- Ein offizieller ISI Standard konnte nicht gefunden werden - vorhanden: Web of Science Core Collection Field Tags von Clarivate
- Eine ISI Tag-Erweiterung über das Format von Clarivate hinaus unter gemeinsamer Spezifikation mit openalexR und bibliometrix Entwicklerinnen und Entwicklern scheint weitere Möglichkeiten zu erschließen.
- Für neue Mappings wären auch in bibliometrix neue Auswertungsalgorithmen und –szenarien zu konzipieren
  
- **Nutzen DeepGreen: Chance zur Nutzung einer noch nie dagewesenen Datenlage durch neue Werkzeuge und Schnittstellen – Reichweitenmessung außerhalb von Verlagen und Contentprovidern durch verteilte Repositorylandschaft – Bildung eigener Metriken auf Artikelebene und teilw. Lösung von kommerzialisierten Metriken. „Open Access auch für Datenanalysen“ Metriken, Sammlungsanalysen, Reichweitenmessung**


Vielen Dank für Ihre Aufmerksamkeit

Vielen Dank für die Unterstützung an die Kollegen der UB Kiel

Dr. Kai Lohsträter  
Adrian Sommer  
Torsten Krause  
Oliver Weiner

Thorsten Wetzenstein  
Universitätsbibliothek Kiel Referat Open-Access-Publizieren | Universitätsverlag |  
Leibnizstr. 9  
24118 Kiel

Telefon: +49 431 880-2740  
mail: [wetzenstein@ub.uni-kiel.de](mailto:wetzenstein@ub.uni-kiel.de)

 <https://orcid.org/0000-0001-9589-3136>

Ursprünglicher Artikel  
Wetzenstein, Thorsten: 1000ste „DeepGreen“-Veröffentlichung an der CAU – ein Zwischenbericht. Feb. 2023

<https://oa-info.sh/2023/02/1000ste-deepgreen-veroeffentlichung-an-der-cau-ein-zwischenbericht/>