

Validating and enriching ^{VABB-SHW} data through external sources

Raf Guns <raf.guns@uantwerpen.be>

ENRESSH training school, Poznań, 24 Oct 2019



Validating enriching
external sources

Manual

19. Complete missing data and validate the accuracy of metadata

“If the use of the actual research output is not feasible, complete missing metadata and validate the accuracy using other national and international bibliographic sources (e.g., CrossRef, WorldCat, Web of Science). Ideally, external data sources should be chosen on the basis of transparency and reliability of their data collection practice.”



Manual



30. Follow and adapt to developments in research practices, research policy, and database maintenance

“Since the context of databases is subject to change, databases need to adapt to ensure that they can continue to fulfil their purposes.”

Example

VABB-SHW: Comprehensive database of SSH publications, Flanders (Belgium)

- Created in context of Flemish funding model to cover SSH, as a complement to Web of Science
- 2019: funding system reform
 - “The international collaboration parameter C3 calculates the share of each university in the number of publications of which one or more addresses on it [...] belong to the university, together with at least one foreign address [...]”

Exercise

What kind of data would, in your opinion, be especially interesting/useful to enrich your database* with?

- Feel free to be creative! No wrong answers 😊
- Take a minute to think and write down some ideas
- Then discuss with your neighbours:
 - Do you have similar ideas? If not, would you add your neighbour's ideas to your list?
 - Do you think these ideas are practically feasible? Why (not)?

* 'your database' = either the one you work on/with in your daily practice, or the one your group has chosen as the focus

How?

Persistent identifiers are key



“For example, if each periodical would assign **unique code numbers** to the articles published, it would be possible for authors to list these numbers in their bibliographies and, thus, to save the work of coding on the part of the citation index staff. It is unlikely that such a development could take place in less than 5 or 10 years, but it is comparable to the problem of getting publishers to include Library of Congress card numbers in their publications.”

– Eugene Garfield, 1955



ORCID



GRID 



DOI

Identifier for (electronic version of) publication or other object

Controlled by a DOI registration agency, the oldest and largest of which is **CrossRef**

Who or what *is* CrossRef anyway?

- Started as an initiative of 12 journal publishers in 2000
- Set up as a non-profit organization, with membership of publishers and other DOI-using organizations
- Each member can continuously update the metadata associated with their DOIs
- Currently ~11,000 members, ~109 million objects

Record

Identificatie c:vabb:322868
Bookmark <https://anet.be/record/opacvabbg/c:vabb:322868/N>

Publicatie

Titel Economic hardship and depression across the life course : the impact of welfare state regimes

Auteur Levecque, Katia
Van Rossem, Ronan
De Boyser, Katrien
Van de Velde, Sarah
Bracke, Piet

Nummer Handle <http://hdl.handle.net/1854/LU-1258270>
doi [10.1177/0022146510394861](https://doi.org/10.1177/0022146510394861)
ISSN 0022-1465
ISI 000291611900008

Uitgave 2011

Omvang 52:2, p. 262-276

Annotatie(s) Verkorte titel tijdschrift: J. Health Soc. Behav.
published

In tijdschrift: Journal of health and social behavior. - Los Angeles. -
J. HEALTH SOC. BEHAV.. - . -

E-info <https://doi.org/10.1177/0022146510394861>

Onderwerp

Collecties VABB. Psychologie
VABB. Sociologie
VABB. Communicatiewetenschappen

<https://api.crossref.org/works/10.1017/s0016774600000688>

```
{
  "status": "ok",
  "message-type": "work",
  "message-version": "1.0.0",
  "message": {
    "indexed": {
      "date-parts": [[2019, 9, 11]],
      "date-time": "2019-09-11T12:00:00Z",
      "timestamp": 1395619200000,
      "delay-in-days": 1209,
      "content-version": "unspecified",
      "content-domain": {
        "domain": "journals.cambridge.org",
        "published-print": {
          "date-parts": [[2010, 12]]
        },
        "abstract": "<jats:title>Abstract</jats:title><jats:p>At the local level, the Maldegem-Stekene coversand ridge and dated using optically stimulated luminescence (OSL). The study aimed at contributing to the understanding of archaeological settlements in this area. The core comprised a 7 m thick series of laminated and massive aeolian sands, in which samples were collected for quartz-based SAR-OSL dating; an internally consistent dataset was obtained. The ages of the lowermost 1 m of the sediments may represent the time-equivalent deposit of a deflation phase that occurred during the Late Pleniglacial and led to the formation of the Gravel Bed. However, a significant part of the sediments (at least 4 m) was deposited later, i.e. during the Aller\u00f8d and the Younger Dryas coversand ridge at the Heidebos locality on the basis of direct age information. The relatively high sedimentation rate and the high aeolian activity and landscape instability during the Late Glacial, which provides part of the environmental framework for human activities in the region (Belgium).</jats:p>",
          "DOI": "10.1017/s0016774600000688",
          "type": "journal-article",
          "created": {
            "date-parts": [[2016, 7, 26]],
            "date-time": "2016-07-26T12:00:00Z",
            "timestamp": 1388520000000,
            "source": "Crossref",
            "is-referenced-by-count": 11,
            "title": ["The timing of aeolian events near archaeological settlements at Maldegem-Stekene (Belgium)"],
            "prefix": "10.1017",
            "volume": "89",
            "author": [
              {
                "given": "C.",
                "family": "Derese",
                "sequence": "first",
                "affiliation": []
              },
              {
                "given": "A.",
                "family": "Zwertvaegher",
                "sequence": "additional",
                "affiliation": []
              },
              {
                "given": "M.",
                "family": "Court-Picon",
                "sequence": "additional",
                "affiliation": []
              },
              {
                "given": "P.",
                "family": "Cromb\u00e9",
                "sequence": "additional",
                "affiliation": []
              },
              {
                "given": "J.",
                "family": "Vermeir",
                "sequence": "additional",
                "affiliation": []
              },
              {
                "given": "A.",
                "family": "haute",
                "sequence": "additional",
                "affiliation": []
              }
            ],
            "member": "56",
            "published-online": {
              "date-parts": [[2016, 7, 26]]
            },
            "asserted-by": "publisher",
            "key": "S0016774600000688_ref021",
            "DOI": "10.1016/S1350-4487(02)00037-9",
            "doi-asserted-by": "publisher",
            "key": "S0016774600000688_ref009",
            "DOI": "10.1017/S0033822200033865",
            "key": "S0016774600000688_ref044",
            "author": "Tavernier",
            "first-page": "159",
            "key": "S0016774600000688_ref018",
            "first-page": "343",
            "article-title": "De morfologie van de Maldegem-Stekene",
            "key": "S0016774600000688_ref022",
            "first-page": "217",
            "article-title": "Bijdrage tot de geologie van de Maldegem-Stekene (Belgi\u00eb)",
            "volume": "40",
            "author": "Heyse",
            "year": "1979",
            "journal-title": "Verhandelingen van de Koninklijke Academie van Wetenschappen (Afdeling Natuurkunde)",
            "key": "S0016774600000688_ref053",
            "doi-asserted-by": "crossref",
            "first-page": "1",
            "DOI": "10.1017/S0016774600000688_ref012",
            "author": "Cromb\u00e9",
            "journal-title": "d'Anthropologie et de Pr\u00e9histoire",
            "key": "S0016774600000688_ref048",
            "DOI": "10.1002/journal-title": "Geoarchaeological research of the large palaeolake of the Moervaart (municipalities of Maldegem and Stekene)",
            "key": "S0016774600000688_ref037",
            "DOI": "10.1016/S1350-4487(02)00037-9",
            "doi-asserted-by": "publisher",
            "key": "S0016774600000688_ref029",
            "DOI": "10.1002/journal-title": "Notae Praehistoricae",
            "key": "S0016774600000688_ref038",
            "DOI": "10.1017/S0033822200034202",
            "doi-asserted-by": "publisher",
            "page": "165",
            "year": "2002",
            "volume-title": "Jutland Archaeological Society Publications",
            "key": "S0016774600000688_ref049",
            "DOI": "10.1016/j.quaint.2007.11.017",
            "journal-title": "Atlas de Belgique",
            "key": "S0016774600000688_ref043",
            "journal-title": "Bulletin de la Soci\u00e9t\u00e9 belge de G\u00e9ologie",
            "volume": "55",
            "author": "Tavernier",
            "year": "1946"
          }
        }
      }
    }
  }
}
```

Uhmmm...

<https://api.crossref.org/works/10.1017/s0016774600000688>

```
▶ update-policy: "http://dx.doi.org/10.117...e-journals-update-policy"
  source: "Crossref"
  is-referenced-by-count: 36
▼ title:
  ▼ 0: "Economic Hardship and Depression across the Life Course"
  prefix: "10.1177"
  volume: "52"
▼ author:
  ▼ 0:
    given: "Katia"
    family: "Levecque"
    sequence: "first"
  ▼ affiliation:
    ▼ 0:
      name: "Research Foundation-Flanders, Brussels, Belg
    ▼ 1:
      name: "Ghent University, Ghent, Belgium"
  ▼ 1:
    given: "Bonen"
```

Tip!

Firefox shows JSON in a human-readable way

<https://api.crossref.org/works/10.1017/s0016774600000688>

```
▼ reference:  
  ▼ 0:  
    key: "bibr1-0022146510394861"  
    DOI: "10.1177/0952872002012002114"  
    doi-asserted-by: "publisher"  
  ▼ 1:
```

How many citations are open today?

Initiative for Open Citations (I4OC)



```
  year: 1998  
  volume-title: "Poverty in Europe"  
  ▼ 3:  
    key: "bibr4-0022146510394861"  
    DOI: "10.1177/1350506809341512"  
    doi-asserted-by: "publisher"
```

<https://api.crossref.org/works/10.1017/s0016774600000688>

```
▼ link:  
  ▼ 0:  
    ▼ URL: "http://journals.sagepub.com/doi/pdf/10.1177/0022146510394861"  
      content-type: "application/pdf"  
      content-version: "vor"  
      intended-application: "text-mining"  
  ▼ 1:  
    ▼ URL: "http://journals.sagepub.com/doi/full-xml/10.1177/0022146510394861"  
      content-type: "application/xml"  
      content-version: "vor"  
      intended-application: "text-mining"  
  ▼ 2:  
    ▼ URL: "http://journals.sagepub.com/doi/pdf/10.1177/0022146510394861"  
      content-type: "unspecified"  
      content-version: "vor"  
      intended-application: "similarity-checking"
```

<https://api.unpaywall.org/v2/10.1017/s0016774600000688?email=your@email.com>

```
best_oa_location:    null
data_standard:      2
doi:                 "10.1177/0022146510394861"
doi_url:             "https://doi.org/10.1177/0022146510394861"
genre:               "journal-article"
has_repository_copy: false
is_oa:               false
journal_is_in_doaj:  false
journal_is_oa:       false
journal_issn_l:      "0022-1465"
journal_issns:       "0022-1465,2150-6000"
journal_name:        "Journal of Health and Social Behavior"
oa_locations:        []
oa_status:           "closed"
```


<https://api.altmetric.com/v1/doi/10.1017/s0016774600000688>

altmetric_id:	582088
schema:	"1.5.4"
is_oa:	false
cited_by_posts_count:	3
cited_by_tweeters_count:	3
cited_by_accounts_count:	3
last_updated:	1421681839
score:	1.75
▼ history:	
1y:	0
6m:	0
3m:	0
1m:	0
1w:	0

EASY

MEDIUM

HARD

Publicatie

Titel The timing of aeolian events near archaeological settlements around Heidebos (Moervaart area, N Belgium)

Auteur Derese, Cilia
Vandenberghe, Dimitri
Zwertvaegher, Ann
Court-Picon, Mona
Verniers, Jacques
Van den haute, Peter

Nummer Handle <http://hdl.handle.net/1854/LU-1268703>
ISI 000293365300001
ISSN 0016-7746

Uitgave 2010

Omvang 89:3-4, p. 173-186

Annotatie(s) Verkorte titel tijdschrift: Neth. J. Geosci.
published

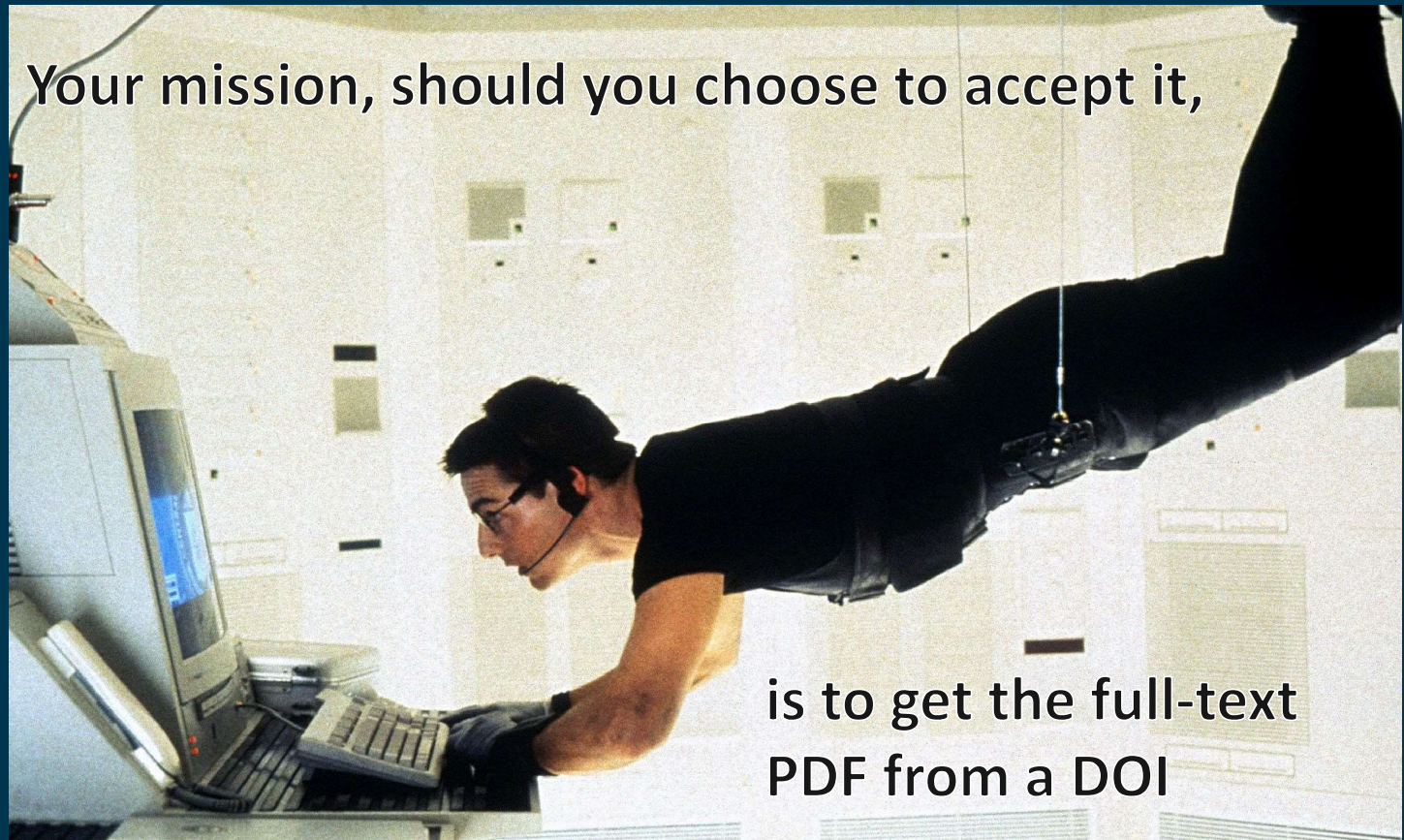
In tijdschrift: *Geologie & mijnbouw.* - Cambridge, [UK]. -

Korte inhoud

At the locality of Heidebos (Moervaart area, N Flanders, Belgium), a sedimentary core was taken in the Malde stimulated luminescence (OSL). The study aimed at contributing to an improved understanding of the evolution in this area. The core comprised a 7 m thick series of laminated and massive aeolian sands, in which several samples were collected for quartz-based SAR-OSL dating; an internally consistent dataset was obtained. The

No DOI! Options?

The messy practice



The messy practice

<https://doi.org/yourdoi> redirects to another document. This may be...

- the full-text PDF 🎉
- a web page that may, somewhere, contain a link to the PDF
 - `<meta name="citation_pdf_url" value="link_to_PDF">` (Google Scholar)
 - Or try all links and see if anything PDF-like comes up...
- CrossRef metadata may contain URLs for text-mining (may require extra steps)
- Some publishers have a fairly straightforward URL scheme (e.g., <https://www.jstor.org/stable/pdf/yourdoi.pdf>)
- Or you can try, e.g., Unpaywall to find an open access PDF

The messy practice

You may also get:

- Other file formats than PDF: XML, e-books, images, plain text, and simply unreadable files...
- SSL errors
- empty pages
- timeouts
- unreachable servers
- blocked

Dear Raf,

Thank you for your email.

I have unblocked the IP address. As you had kindly pointed out, you had triggered this IP block by multiple actions from an IP in a very short period of time. This type of usage can indicate illegal activity and may be a violation. It is for that reason that we place a block against it.

If you have any further questions, please do not hesitate to contact us.

Kind regards,

[Redacted Signature]

Taylor & Francis Online Customer Services

www.tandfonline.com

Enriching may introduce new problems

Number of publications with structured affiliation data by country of origin of the (co-)authors:

754

NUMBER OF PUBLICATIONS PER COUNTRY OF INSTITUTION OF AUTHORS

Country (48)	Number
Belgium	750
Netherlands	77
United States	22
United Kingdom	19

Enriching may introduce new problems

Alissa De Ceuninck is Master in Film Studies and Visual Culture at the University of Antwerp. Earlier, she obtained a Master in Communication Studies at the Free University of Brussels, specializing in media, democracy and journalism.

M. Morrens is connected as doctor training to be a psychiatrist to the Collaborative Antwerp Psychiatric Research Institute (CAPRI), Campus Drie Eiken, Antwerp.

Are these affiliations?

Why (not)?

Which institute(s) are they affiliated to?

More pitfalls

- External databases may be wrong while your database is right
 - Don't blindly overwrite all values that are different!
- Not everything has a persistent identifier
 - Especially publications in local journals, publications in books or proceedings, gray literature...
 - VIRTIA (Finland): 67% of peer-reviewed outputs has DOI, only 22% for book publications
 - Cuban journals cannot register DOIs because of trade embargo

We have new data. Now what?

Store all data in own database or **fetch** latest data from, e.g., CrossRef when needed?

- Many data are dynamic, prone to change
- But fetching data when needed is only useful for 'one at a time' usage (e.g. showing some extra data in a record)

→ Often storage with regular updates will be needed

We have new data. Now what?

“For databases where data are collected by means of data transfer, it is important to consider the relationship between the data in the national database and in the database from which the data originated. If the data are enriched in the national databases, it is useful to **implement procedures that allow to improve also the accuracy of data in the databases from which the data originated**. This, however, requires coordination between different organisations, consideration of the ownership of data as well as different legal frameworks that might influence this process.”

- Implement feedback loop from national database to source database(s)
 - Not always possible for legal, technical... reasons
 - Extra validation by source organisations may be needed

Take-away messages

- Validating and enriching is useful but not always easy
- Persistent identifiers are core infrastructure to enable this. Objects that lack them are much harder to compare to external systems.
- Don't blindly trust the external system for corrections.
- Be aware of dynamicity of data, i.e., think carefully about storage and 'refreshing'.
- Implement feedback loop to source databases, if needed and possible.