

# Data structure in citation files

Jonathan Hooper, Leeds. May 2022

## Contents

Overview .....	1
Description .....	1
Citation .....	1
Citation.Course .....	2
Citation.Leganto .....	3
Citation.Alma .....	4
Citation.Scopus .....	5
Citation.Scopus.first-match .....	6
Citation.Scopus.first-match.Author .....	6
Citation.WoS element .....	7
Citation.WoS.first-match .....	9
Citation.WoS.first-match.metadata.Author .....	9
Citation.WoS.first-match.metadata.Address .....	11
Citation.WoS.first-match.metadata.Reprint_address .....	11
Citation.VIAF element .....	11
Citation.VIAF.best-match .....	12

## Overview

Citation enhancement consists of running a number of separate scripts, which progressively enhance the original reading list data with data from external sources.

The format used to temporarily save this data in-between steps in the process is JSON.

In general, a JSON file will exist in the directory **Data/** for each module being analysed. This file lists the citations in that module's reading lists (the module may have zero, one or more reading lists).

## Description

The data structure in the JSON citations files is not centrally defined or enforced, and is simply based on what is convenient when extracting data from the various sources. Defining a standard data format, and a standard for the various data integrations, might be a future improvement.

Each JSON file contains an array of citation elements i.e. the file looks like:

```
{ Citation, Citation, Citation, ... }
```

## Citation

Normally a **Citation** element represents an individual citation in a reading list:

- The exceptions to this are:
  - Where we try to collect information for a module that does not have a corresponding course in Alma
  - Where any corresponding course(s) in Alma do not have any associated reading lists
  - Where a associated reading list has zero citations
  - In those cases, a **single** place-holder **Citation** element will be present for the module, without any citation-level data
- A citation may simply be a note, without any bibliographic reference
  - In this case there will be a Citation element for that note, but without any data enhancement from Alma, Scopus, WoS or VIAF
  - *Notes can be distinguished from the value of Citation.Leganto.secondary\_type.value ("NOTE")*

A citation element is an associative array with possible keys **Course**, **Leganto**, **Alma**, **Scopus**, **WoS** and **VIAF** i.e.:

```
{ "Course":Course, "Leganto":Leganto, "Alma":Alma, "Scopus":Scopus, "WoS":WoS, "VIAF":VIAF }
```

- **Course** will always be present
- **Leganto** will *usually* be present (if a reading list with citations exists)
- The other keys will be present once a search has been attempted against that source for this citation

These top-level values in the Citation associative array are described below:

#### Citation.Course

Added to citation by: **getCitationsByModule.php**

Holds information about the **module**

#### Properties:

Property	Type	Present?	Meaning
<b>modcode</b>	String	Always	<p>The short module code (e.g. <b>DEPT1234</b>) this citation is associated with</p> <p><i>One JSON file may in principle contain citations from more than one module code.</i></p> <p><i>But in practice, we tend to keep to one file per module, and the batch script will always use one file per module.</i></p>
<b>course_code</b>	String	If a course exists in Alma for modcode	<p>The course code this citation's reading list is associated with, as used in Alma</p> <p><i>If not present, this means no course exists in Alma for this</i></p>

			<i>modcode - and no further enhancement can or will be done</i>
<b>list_code</b>	String	Always	The code for the associated reading list <i>This duplicates the entry in Citation.Leganto (below) because of the case where a list has zero citations, and so has no Citation.Leganto element</i>

### Citation.Leganto

Added to citation by: **getCitationsByModule.php**

Holds information about an individual **reading list citation**. Will only be present if there is a course in Alma matching the modcode, **and** if that course has a reading list.

### Properties:

Property	Type	Present?	Meaning
<b>id</b>	Integer	Always	The system-generated ID for this reading list
<b>list_code</b>	String	Always	The code for this reading list
<b>list_title</b>	String	Always	The title of the reading list. Unlike id and list_code, this can be changed by instructor/library and displays to users, so this may be a better way of referring to a reading list in user-facing reports
<b>citation</b>	Integer	Always	Position of citation in list (starting at 1)  This should always correspond to ordering of citations in Alma and Leganto, when viewing in instructor-order.  <i>Notes count as citations for the purposes of this numbering.</i>
<b>section</b>	String	Always	The name of the reading list section this citation belongs to  <i>May be helpful in reports, especially for long reading lists</i>
<b>secondary_type</b>	Associative array	Always	The citation's type code and description as used in Alma e.g. <b>{"value": "BK", "desc": "Book"}</b>  <i>Types are listed in Ex Libris's <a href="#">Reading List Loader</a> documentation.</i>
<b>metadata</b>	Associative array	Always	Selected fields from the citation data returned by the Alma Courses API e.g. <b>title, author, mms_id</b>

			<i>getCitationsByModule.php could be modified to fetch additional fields if required</i>
--	--	--	------------------------------------------------------------------------------------------

### Citation.Alma

Added to citation by: **enhanceCitationsFromAlma.php**

Holds information about a **bibliographic record** from Alma. Will only be present if there is a value in **Citation.Leganto.metadata.mms\_id** with a corresponding bibliographic record in Alma.

#### Properties:

Property	Type	Present?	Meaning
<b>titles</b>	Array of associative arrays	Usually	<p>Title-like data from the bibliographic record (<i>defined by Marc tags in enhanceCitationsFromAlma.php</i>).</p> <p>May be more than one value e.g. title plus alternate title. And may be alternative representations (e.g. in Chinese characters) from linked 880 field.</p> <p>For each title:</p> <p><b>title.tag</b> holds the Marc (or linked Marc) tag <i>e.g.</i> "245"  <b>title.originalTag</b> may hold "880" if the data comes from a linked field  <b>title.a</b> holds the value of the \$a subfield if present <i>etc</i>  <b>title.collated</b> holds the values of the \$a \$b subfields, joined with spaces</p>
<b>creators</b>	Array of associative arrays	Usually	<p>Author-like data from the bibliographic record (defined by Marc tags in enhanceCitationsFromAlma.php).</p> <p>May be more than one value e.g. author plus additional author. And may be alternative representations (e.g. in Chinese characters) from linked 880 field.</p> <p>For each creator:</p> <p><b>creator.tag</b> holds the Marc (or linked Marc) tag <i>e.g.</i> "100"  <b>creator.originalTag</b> may hold "880" if the data comes from a linked field  <b>creator.a</b> holds the value of the \$a subfield if present <i>etc</i>  <b>creator.collated</b> holds the values of the \$a \$b \$c \$d \$q subfields, joined with spaces</p>
<b>ids</b>	Associative array	Often	Identifiers (ISBNs, ISSN, LCCNs) found in the bibliographic record

			<i>e.g. <b>ids.isbn</b> holds an array of ISBN values</i>
--	--	--	-----------------------------------------------------------

## Citation.Scopus

Added to citation by: **enhanceCitationsFromScopus.php**

Holds information about **Scopus searches and results** for a citation.

## Properties:

Property	Type	Present?	Meaning
<b>rate-limit</b>	Associative array	Usually	Rate-limit data from the API, with an entry for each type of query attempted with this citation.  For each query-type: <b>limit</b> = maximum number of allowed requests per week <b>remaining</b> = number of allowed requests remaining this week <b>reset</b> = timestamp of start of next week period  <i>NB when any of the <b>remaining</b> figures reach zero, we will no longer be able to harvest data until the start of the next week (i.e. the time in the <b>reset</b> timestamp)</i>
<b>searches-no-results</b>	Array	Yes (may be empty)	Scopus API queries that returned no results for this citation, in the order they were attempted  <i>The script tries progressively broader searches until it gets a result - see 2_2_searching_the_data_sources.pdf</i>
<b>search-active</b>	String	If successful search	The first Scopus API query that retrieved at least one result
<b>search-pref</b>	Integer	If successful search	The position of the first successful query in the order-of-preference – a lower number indicates a narrower search
<b>result-count</b>	Integer	If successful search	The number of records returned by the first successful query  <i>NB currently the code is only analysing the first result - it does not attempt to determine whether a result lower-down the set may be a better match to the reading list citation</i>
<b>results</b>	Array of associative arrays	If successful search	Brief summary of all the records returned by the first successful search, using selected data from the API response

<b>first-match</b>	Associative array	If successful search	Data from the API response for the first record in the result set from the first successful search - this is the record that will be used in subsequent processing.  <i>See Citation.Scopus.first-match description below</i>
--------------------	-------------------	----------------------	-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

### [Citation.Scopus.first-match](#)

Added to citation by: **enhanceCitationsFromScopus.php**

Holds information about the **Scopus record** that will be used in subsequent processing.

#### Properties:

Property	Type	Present?	Meaning
<b>summary</b>	Associative array	Always	Core metadata for this record, selected by the script from the full Scopus record (id, title, creator, publication, doi, document type).
<b>self</b>	String [URL]	Yes	A URL identifying this record in the Scopus API
<b>authors</b>	Array	Always (?)	Data on the authors of this citation i.e. [ <b>Author, Author, Author, ...</b> ]  <i>See Citation.Scopus.first-match.Author description below</i>

### [Citation.Scopus.first-match.Author](#)

Added to citation by: **enhanceCitationsFromScopus.php**

Holds information about an individual **author** from the Scopus record that will be used in subsequent processing. Typically a single citation may contain several authors.

#### Properties:

Property	Type	Present?	Meaning
<b>@auid</b>	String [numeric]	Always	Author ID in Scopus
<b>author-url</b>	String [URL]	Yes	A URL identifying this author in the Scopus API
<b>preferred-name</b>	Associative array	Always (?)	The preferred author-name  <i>The similarity calculation may try various forms of name against the reading list authors, and the export script will decide which form(s) to use in the report</i>
<b>ce:initials, ce:surname, ce:indexed-name</b>	String	Always (?)	Author-name data  <i>The similarity calculation may try various forms of name against the reading list authors, and the export script will decide which form(s) to use in the report</i>

<b>affiliation</b>	Associative array	Often	<p>The author's <b>contemporary</b> affiliation (i.e. at the time of the publication), broken into subfields</p> <p><b>affiliation.country</b> generally contains standard English-language country names, but other forms may appear.</p> <p>No conversion/translation is done in saving this data from Scopus, it is up to the export scripts to handle this.</p>
<b>affiliation-current</b>	Associative array	Often	<p>The author's <b>current</b> affiliation (i.e. at the time of data extract), broken into subfields</p> <p>NB the subfield structure of affiliation and affiliation-current are slightly different, because of the different way the data is received from the API.</p> <p><b>affiliation-current.address.@country</b> generally contains ISO-3-letter country codes, but other forms may appear.</p> <p>No conversion/translation is done in saving this data from Scopus, it is up to the export scripts to handle this.</p>
<b>similarity-title</b>	Integer	Always	<p>Scores 0-100 giving string-similarity between title as appears in reading list and title as appears in Scopus record</p> <p>Value will be identical for all authors in a given citation (because all Scopus authors belong to the same Scopus citation).</p> <p>See <a href="#">2_3_author_title_similarity_scores.pdf</a> for details of Similarity calculation</p>
<b>similarity-author</b>	Integer	Always	<p>Scores 0-100 giving string-similarity between author as appears in reading list and author as appears in Scopus record</p> <p>Unlike similarity-title, will vary between authors in a citation.</p> <p>See <a href="#">2_3_author_title_similarity_scores.pdf</a> for details of Similarity calculation</p>
<b>search-parameters</b>	Associative array	Always	<p>The data values (DOI, TITLE etc) extracted from Leganto/Alma, to be used in searching Scopus</p>
<b>extra-parameters</b>	Associative array	Always	<p>Data values extracted from Leganto/Alma that are <b>not</b> used in searching Scopus, but <i>may</i> be used by the script to calculate similarity</p>

Citation.WoS element

Added to citation by: **enhanceCitationsFromWos.php**

Holds information about WoS searches and results for a citation.

**Properties:**

Property	Type	Present?	Meaning
<b>rate-limit</b>	Associative array	Usually	<p>Rate-limit data from the API.</p> <p><b>x-rec-amtperyear-remaining</b> = searches remaining this year</p> <p><b>x-req-reqpersec-remaining</b> = searches remaining this second</p> <p><i>NB if either figure reaches zero, we will no longer be able to harvest data in that period.</i></p> <p><i>I assume the per-year quota will be reset when the annual subscription is renewed (c. April 2023?)</i></p>
<b>searches-no-results</b>	Array	Yes (may be empty)	<p>Searches using the WoS API Expanded query syntax that returned no results for this citation, in the order they were attempted</p> <p><i>The script tries progressively broader searches until it gets a result - see 2_2_searching_the_data_sources.pdf</i></p>
<b>search-active</b>	String	If successful search	The first search that retrieved at least one result, using the WoS API query syntax
<b>search-pref</b>	Integer	If successful search	The position of the first successful search in the order of preference of search strategies – a lower number indicates a narrower search
<b>result-count</b>	Integer	If successful search	<p>The number of records returned by the first successful search</p> <p><i>NB currently the code is only analysing the first result - it does not attempt to determine whether a result lower-down the set may be a better match to the reading list citation</i></p>
<b>results</b>	Array of associative arrays	If successful search	Brief summary of all the records returned by the first successful search, using selected data from the API response
<b>first-match</b>	Associative array	If successful search	<p>Data from the API response for the first record in the result set from the first successful search - this is the record that will be used in subsequent processing.</p> <p><i>See Citation.WoS.first-match description below</i></p>



[Citation.WoS.first-match](#)

Added to citation by: **enhanceCitationsFromWos.php**

Holds information about the WoS record that will be used in subsequent processing.

**Properties:**

Property	Type	Present?	Meaning
<b>metadata</b>	Associative array	Always	<p>Core metadata for this record, selected by the script from the full WoS record (id, source, title, year, doctype, publisher, authors, addresses etc).</p> <p>Particular individual entries to note:</p> <ul style="list-style-type: none"> <li>• <b>authors</b> is an array of individual authors for this citation i.e. [ <b>Author, Author, Author...</b> ] <i>See Citation.WoS.first-match.metadata.Author description below</i></li> <li>• <b>addresses</b> is an array of individual addresses for this citation i.e. [ <b>Address, Address, Address...</b> ] <i>See Citation.WoS.first-match.metadata.Address description below</i></li> <li>• <b>reprint_addresses</b> is an array of individual reprint (corresponding) addresses for this citation i.e. [ <b>Reprint_address, Reprint_address, Reprint_address...</b> ] <i>See Citation.WoS.first-match.metadata.Reprint_address description below</i></li> </ul>

[Citation.WoS.first-match.metadata.Author](#)

Added to citation by: **enhanceCitationsFromWos.php**

Holds information about an individual author from the WoS record that will be used in subsequent processing. Typically a single citation may contain several authors. Where an author is linked to a specific address (generally post-2008) also includes author's address - otherwise all addresses are in CITATION.WoS.first-match.metadata.addresses (see below).

**Relevant properties:**

Property	Type	Present?	Meaning
<b>addr_no</b>	Integer   String [space-separated Integers]	Often (post-2008)	<p>Identifies the address(es) in CITATION.WoS.first-match.metadata.addresses for this author</p> <p><i>If Integer identifies a single address e.g. 12 -&gt; Address # 12</i></p>

			<i>If String list of Integers identifies a list of addresses e.g. e.g. "2 3 5" -&gt; Addresses # 2, 3, 5</i>
<b>last_name, display_name, full_name, wos_standard, first_name</b>	Strings	Usually	Various forms of author's name  <i>The similarity calculation may try various of these against the reading list authors, and the export script will decide which form(s) to use in the report</i>
<b>preferred_name</b>	Associative array	Always (?)	full_name, last_name, middle_name, first_name  <i>The similarity calculation may try various of these against the reading list authors, and the export script will decide which form(s) to use in the report</i>
<b>addresses</b>	Array of associative arrays	Often (post-2008)	Copy (for convenience) of the address from Citation.WoS.first-match.metadata.addresses which matches this author (where available).  Array of addresses, each one an associative array of different elements  <i>address.country generally contains standard English-language country names but other forms may appear - no conversion/translation is done in saving this data from WoS, it is up to the export scripts to handle this</i>
<b>similarity-title</b>	Integer	Always	Scores 0-100 giving string-similarity between title as appears in reading list and title as appears in WoS record  Value will be identical for all authors in a given citation (because all WoS authors belong to the same Scopus citation).  <i>See 2_3_author_title_similarity_scores.pdf for details of Similarity calculation</i>
<b>similarity-author</b>	Integer	Always	Scores 0-100 giving string-similarity between author as appears in reading list and author as appears in WoS record  Unlike similarity-title, will vary between authors in a citation.  <i>See 2_3_author_title_similarity_scores.pdf for details of Similarity calculation</i>

### Citation.WoS.first-match.metadata.Address

Added to citation by: **enhanceCitationsFromWos.php**

Holds information about an individual address from the WoS record. Typically a single citation may contain several addresses: Citation.WoS.first-match.metadata.Address always contains **all** addresses in the citation, whether or not individual authors are linked to them.

#### Properties:

Property	Type	Present?	Meaning
<b>names</b>	Array of associative arrays	Usually	Authors associated with this address  <i>NB We do not refer to this data and instead come in the other direction, from the authors array, and from there to linked addresses</i>
<b>address_spec</b>	Associative array	Always	Details of this address (including country)  <i>NB address_spec.addr_no forms the link with Author.addr_no.</i>  <i>There may be a many-to-many mapping between authors and addresses. There may also be addresses that are not linked to any individual author.</i>

### Citation.WoS.first-match.metadata.Reprint\_address

Added to citation by: **enhanceCitationsFromWos.php**

Holds information about an individual reprint address (corresponding address) from the WoS record.

*We are currently not making use of these addresses in the reports we generate, because we understand they do not add to the information in Addresses (unless the reprint address is a purely administrative one).*

#### Properties:

Property	Type	Present?	Meaning
<b>names</b>	Array of associative arrays	Usually	Authors associated with this address
<b>address_spec</b>	Associative array	Always	Details of this address (including country)

### Citation.VIAF element

Added to citation by: **enhanceCitationsFromViaf.php**

Holds information about VIAF searches and results for a citation.

#### **Important differences between VIAF, and Scopus/WoS results**

- Scopus and WoS are **citation**-databases where VIAF is a **person**-database:
- With Scopus and WoS, a **single** search is done at the **citation-level**
- With VIAF, a **separate** search is done **for each author** in the reading list citation

- *Citation.VIAF is an **array** of datasets, each one for a separate search for a different author*
- *Citation.Scopus and Citation.WoS are each single datasets for a citation - but each of these single datasets may contain data for multiple authors*
- *Different VIAF authors for a citation will typically have **different title-similarity scores***
- *Different Scopus and WoS authors will have the **same** title-similarity*
- *With Scopus and WoS, the code currently just takes the first record from a result set - with VIAF, the code finds the record which best matches by title the reading list citation*

**Properties of each individual entry in the VIAF element:**

Property	Type	Present?	Meaning
<b>searches</b>	Associative array	Always	How the VIAF search is constructed - author and title will be filled in from search-term-au and search-term-ti ( <i>see below</i> )
<b>search-pref</b>	Integer	If successful search	The position of the first successful search in the order of preference of search strategies – a lower number indicates a narrower search <i>e.g. a search-relation "any" in searches.AU or searches.TI is a <b>broad</b>er search than a search-relation "all", and will have a <b>higher</b> search-pref</i>
<b>data-source</b>	String	Always	Where the author/title comes from to carry out the search against VIAF <i>i.e. Alma   Leganto   Scopus   WoS</i>
<b>records</b>	Integer	If successful search	The number of records returned by the first successful search  <i>NB For VIAF, unlike Scopus and WoS, the code looks at all the returned records and decides which is the best match on the basis of title-similarity</i>
<b>results</b>	Array of associative arrays	If successful search	List of author names returned by the first successful search, using selected data from the API response  <i>Only one of these will be the author selected by the script, on the basis of title-similarity</i>
<b>best-match</b>	Associative array	If successful search	Data from the API response for the record in the result set from the first successful search which best matches (by title) the reading list citation. This is the record that will be used in subsequent processing.  <i>See <b>Citation.VIAF.best-match</b> description below</i>

[Citation.VIAF.best-match](#)

Added to citation by: **enhanceCitationsFromViaf.php**

Holds information about a VIAF record that will be used in subsequent processing.

**Properties:**

Property	Type	Present?	Meaning
<b>type</b>	String	Always	The integration script filters this to only include entries with value "Personal"  <i>Some searches will include results for institutions etc, and these need filtering out</i>
<b>heading</b>	String	Always	A definitive form of the personal name  <i>The script prefers the form from Library of Congress, otherwise it takes the first one found</i>
<b>headings-all</b>	Array of Strings	Always	All the forms of the personal name present in the record  <i>The \$a \$b \$d \$q sub-fields of the VIAF data, joined with spaces</i>
<b>about</b>	String [URL]	Always	A URL identifying this record in the VIAF API
<b>similarity-title</b>	Integer	Always	Scores 0-100 giving string-similarity between title as appears in reading list, and a title associated with this author in VIAF  <i>VIAF lists works associated with an author, and we compare each entry in this list with each form of the title in Leganto/Alma</i>  <i>See 2_3_author_title_similarity_scores.pdf for details of Similarity calculation</i>
<b>similarity-author</b>	Integer	Always	Scores 0-100 giving string-similarity between author as appears in reading list and author as appears in VIAF record  <i>See 2_3_author_title_similarity_scores.pdf for details of Similarity calculation</i>
<b>best-matching-title</b>	String	Usually	The title associated with this author in VIAF that is the most similar to a title in our reading list
<b>titles</b>	Array of Strings	Usually	All the titles associated with this author in VIAF
<b>nationalities</b>	Array	Often	Nationalities of this author

			<p><i>Generally 2-letter ISO codes (e.g. "GB") but other forms appear, including foreign-language country names</i></p> <p><i>No conversion/translation is done in saving this data from VIAF, it is up to the export scripts to handle this</i></p>
<b>countriesOfPublication</b>	Array	Often	<p>Countries of publication for this author</p> <p><i>Not currently used by export scripts</i></p>
<b>affiliations</b>	Array	Sometimes	<p>Author affiliations (taken from 551\$a fields in VIAF)</p> <p><i>Does not include countries, and so currently not used by export scripts - a future improvement to the process would be to make use of this data</i></p>
<b>locations</b>	Array	Sometimes	<p>Author locations (taken from 510.2_\$a fields in VIAF)</p> <p><i>Does not include countries and so currently not used by export scripts - a future improvement to the process would be to make use of this data</i></p>