# Scraping-Fun

April 19, 2020

## 1  A Scraper for the library records of HAB Wobü

This little program will extract the metadata from the library catalog in structured form. This metadata is supposed to be reused, yet copy and paste is tedious. This script is supposed to automate the process.

### 1.1  Description of the site

- The links looks as follows: http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST=31
- with the =31 at the end being the results with which it starts, there are 71 results linked to Michael Maier.
- possibly not all are needed
- annoyingly, the URL is not really straightforward: http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=37/MAT=/NOMAT=T/REL?PPN=080093043
- Maier, Michael has a fixed person id apparently: `Maier, Michael, 1568 - 1622 (Zeit, Lebensdaten)` = PPN=080093043

```
[1]: from lxml import html
     import requests
```

```
[2]: starting_url = 'http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?
      ↪FRST=01'
     url_front = 'http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST='

     pages_list = ['01','11','21','31','41','51', '61']

     for item in pages_list:
         print(url_front + item)

     page_store = []

     the_pages = [(url_front + item) for item in pages_list]
     the_pages
```

```
http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST=01
http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST=11
http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST=21
```

```
http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST=31
http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST=41
http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST=51
http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST=61
```

[2]:
```
['http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST=01',
 'http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST=11',
 'http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST=21',
 'http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST=31',
 'http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST=41',
 'http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST=51',
 'http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=21/NXT?FRST=61']
```

[3]:
```python
def print_elem(elem):
    print("<%s>\nTags: %s;\n%s...\n\n" % (elem.tag, elem.attrib, elem.
    text_content()[:200].replace('\n', ' ')))
```

[4]:
```python
for one_page in the_pages:
    page = requests.get(one_page)
    page_store.append(page)
```

[5]:
```python
page_store
for page in page_store:
    tree = html.fromstring(page.text)
    print_elem(tree)
```

```
<html>
Tags: {};
        window.addEventListener("load", function(){
window.cookieconsent.initialise({   "palette": {     "popup": {
"background": "#dbdbdb"      },      "button": {       "background": "#9e9e9e"
…
<html>
Tags: {};
        window.addEventListener("load", function(){
window.cookieconsent.initialise({   "palette": {     "popup": {
"background": "#dbdbdb"      },      "button": {       "background": "#9e9e9e"
…
<html>
Tags: {};
        window.addEventListener("load", function(){
window.cookieconsent.initialise({   "palette": {     "popup": {
"background": "#dbdbdb"      },      "button": {       "background": "#9e9e9e"
…
<html>
Tags: {};
        window.addEventListener("load", function(){
window.cookieconsent.initialise({   "palette": {     "popup": {
```

```
"background": "#dbdbdb"        },        "button": {          "background": "#9e9e9e"
…
<html>
Tags: {};
        window.addEventListener("load", function(){
window.cookieconsent.initialise({   "palette": {        "popup": {
"background": "#dbdbdb"        },        "button": {          "background": "#9e9e9e"
…
<html>
Tags: {};
        window.addEventListener("load", function(){
window.cookieconsent.initialise({   "palette": {        "popup": {
"background": "#dbdbdb"        },        "button": {          "background": "#9e9e9e"
…
<html>
Tags: {};
        window.addEventListener("load", function(){
window.cookieconsent.initialise({   "palette": {        "popup": {
"background": "#dbdbdb"        },        "button": {          "background": "#9e9e9e"
…
```

## 2  Overview page of the search results

```python
mainpage = 'http://opac.lbs-braunschweig.gbv.de/DB=2/SET=2/TTL=37/MAT=/NOMAT=T/
  ↪REL?PPN=080093043'
page = requests.get(mainpage)
tree = html.fromstring(page.text)
print_elem(tree)
tree.text_content()
```

```
<html>
Tags: {};
        window.addEventListener("load", function(){
window.cookieconsent.initialise({   "palette": {        "popup": {
"background": "#dbdbdb"        },        "button": {          "background": "#9e9e9e"
…
```

[6]: '\n\n\n\n\n\n\nwindow.addEventListener("load", function(){
window.cookieconsent.initialise({\n   "palette": {\n      "popup": {\n
"background": "#dbdbdb"\n      },\n      "button": {\n        "background":
"#9e9e9e"\n      }\n   },\n   "theme": "classic",\n   "content": {\n      "message":
\'\',\n     "dismiss": \'\',\n     "link": \'\',\n     "href": \'\'\n
}\n})});\n\n\n\n    OPC4 - results/shortlist\n    \n    \n    \n    \n    \n
\n    \n\n\n\n\n\n\n    \n\n    \n\t\n\t\n\t\n\t\n    \n\n    \n

```

```
\n\n\n\xa0\nvar ah=screen.availHeight-100;\nvar
aw=screen.availWidth-50;\ndocument.write(\'<a class="nav0"
href="START_WELCOME">Suchen</a>\');\n\n\n
Suchen\n\n\n\xa0\n\xa0|\xa0\n\xa0Suchergebnis\xa0\n\xa0|\xa0\n\xa0\nvar
ah=screen.availHeight-100;\nvar aw=screen.availWidth-50;\ndocument.write(\'<a
class="nav0" href="ADVANCED_SEARCHFILTER">Erweiterte Suche</a>\');\n\n\n
Erweiterte Suche\n\n\n\xa0\n\xa0|\xa0\n\xa0\nvar ah=screen.availHeight-100;\nvar
aw=screen.availWidth-50;\ndocument.write(\'<a class="nav0"
href="ERROR_MYSHELF_EMPTY">Zwischenablage</a>\');\n\n\n
Zwischenablage\n\n\n\xa0\n\xa0|\xa0\n\xa0\nvar ah=screen.availHeight-100;\nvar
aw=screen.availWidth-50;\ndocument.write(\'<a class="nav0"
href="javascript:void(0)" onclick="javascript:window.open(\\\'/loan/DB=2/LNG=DU/
LRSET=1/MAT=/SET=1/SID=57f7b86d-0/TTL=1/LNG=DU/USERINFO_LOGIN\\\',\\\'loan\\\',\
\\'scrollbars=yes,top=10,left=10,height=\'+ah+\',width=\'+aw+\',resizable=yes,to
olbar=no,location=no,status=no,menubar=no\\\')">Benutzerkonto</a>\');\n\n\n
Benutzerkonto\n\n\n\xa0\n\xa0|\xa0\n\xa0\nvar ah=screen.availHeight-100;\nvar
aw=screen.availWidth-50;\ndocument.write(\'<a class="nav0"
href="HELP_SEARCH">Hilfe</a>\');\n\n\n    Hilfe\n\n\n\xa0\n\n    ©\n     \n
\n\n\n    \n\n\n\n\n\n    \n    \n    \n\n\n    \n\n    \n\t    suchen
[oder]\n\t    suchen [und]\n\t
\n\t\teingrenzen\n\t\terweitern\n\t\tausgenommen\n\t\tneu ordnen\n\t    \n\t
Index blättern\n\t    \n\t\n\t\n\t    \n\t\t\n\t\t[ALL] Alle Wörter\n\t
\n\t\t\n\t\t[PER] Person\n\t    \n\t\t\n\t\t[TIT] Titel (Stichwort)\n\t
\n\t\t\n\t\t[WTP] Werktitel (Phrase)\n\t    \n\t\t\n\t\t[SER] Serie, Zeitschrift
(Stichwort)\n\t    \n\t\t\n\t\t[KOR] Körperschaft, Konferenz, Geografikum
(Stichwort)\n\t    \n\t\t\n\t\t[NUM] Nummern (allgemein)\n\t
\n\t\t\n\t\t[SLW] Schlagwörter\n\t    \n\t\t\n\t\t[BKL] Basisklassifikation\n\t
\n\t\t\n\t\t[SGN] Signatur ohne Blanks und ohne Sonderzeichen\n\t
\n\t\t\n\t\t[VER] Veröffentlichungsangaben\n\t    \n\t\t\n\t\t[BBG] Bibliogr.
Gattung und Status\n\t    \n\t\t\n\t\t[FPR] Fingerprint (Phrase)\n\t
\n\t\t\n\t\t[PRN] Provenienz (Ex-Ebene)\n\t    \n\t\t\n\t\tSystematik Altbestand
[LSY]\n\t    \n              \n\n\npu="";\npopup="";\nfunction PU(URL,w,h,NAME)
{\n    var v=parseInt(navigator.appVersion);\n    if (v>="3") {\n        if
(typeof NAME == \'undefined\' || NAME == "") {\n            NAME="_blank";\n
}\n        WINDOW=\'top=10,left=10,width=\'+w+\',height=\'+h+\',scrollbars=1,res
izable=1,status=0,toolbar=0\';\n        popup=window.open(URL,NAME,WINDOW);\n
popup.focus();\n        if (pu!=popup) {\n            pu=popup;\n        }\n
} else {\n        window.location.href=URL;\n    }\n}\nfunction gbvhelpselect()
{\n  var x = document.SearchForm.IKT.selectedIndex;\n  var help_docs = new
Array(\n        "hab/gbv-1016.html","hab/gbv-1004.html","hab/gbv-4.html","hab/gb
v-8063.html","hab/gbv-5.html","hab/gbv-1005.html","hab/gbv-1007.html","hab/gbv-5
040.html","hab/gbv-5004.html","hab/gbv-2140.html","hab/gbv-2100.html","hab/gbv-8
600.html","hab/gbv-8092.html","hab/gbv-3080.html","hab/gbv-20.html","gbv-9999.ht
ml");\n  var help_url = \'http://opac.lbs-
braunschweig.gbv.de:80/hermes/gbvhelp/du/\'+help_docs[x];\n  var w=window.open(h
elp_url,\'hilfe\',\'scrollbars=yes,height=500,width=800,resizable=yes,toolbar=no
,location=no,status=no,menubar=no\');\n}\ndocument.write(\'<a
```

href="javascript:gbvhelpselect();"><IMG name=gbvhelp src="http://opac.lbs-braunschweig.gbv.de:80/img_psi/2.0/icons/icon_gbvhelp.gif" hspace=0 vspace=0 border=0 alt="Hilfe zum gew&auml;hlten Suchschl&uuml;ssel"></a> \');\n\n\n
\n    \n\n    \xa0sortiert nach\n\t\n\t\n\t    \n\t    Erscheinungsjahr\n\t
\n\t    \n\t    Relevanz\n\t    \n\t\n\n\n\n    \n\n\n\n    \n\t\n\t\n\t\n\t\n
\n    \n\n\n\n  \n\n\n\n    \n\n\n\n\n    \n    \n    \n\n\n    \n    \n\n
Suchgeschichte\n\n    \n         \n    \n\n    Kurzliste\n\n    \n        \n
\n\n    Titeldaten\n\n    \n    \n\n\n    \n\n    \n    \n    \n\n    \n\n    \n
\n\n\n    \n\n\n\n\n    \n\n\n  Katalogmenü\n  \nSpeichern\n
\ndocument.write(\'<TR><TD class="mnu"><p class="mnu"><a class="mnu"
href="EXIT?DEST=https%3A%2F%2Fopac.lbs-braunschweig.gbv.de%2FDB%3D2%2F">Abmelden
</A></TD></TR>\');\n\n\n\nAbmelden\n\nTrefferanalyse\n\n\xa0\n\xa0\n\xa0\n\xa0\n
\xa0\n\xa0\n\xa0\n\xa0\n\n   \n    \n    \n\n\n    \n    \n\n
\n\n1\xa0-\xa010\xa0von\xa069\xa0\xa0\xa0\xa0\xa0\xa0\n    \n\n     Ihre Aktion\n
\xa0bezogen auf\xa0Maier, Michael, 1568 - 1622 (Zeit, Lebensdaten)\n\n\n\n\n
\n    \n    \n\n\n\n    \n    \n\n\n\n    \n    1.\xa0\n    \n\nEin
religionswissenschaftlicher Kommentar zu den Arcana Arcanissima und der
Mythoalchemie des alchemo-hermetischen Iatrochemikers Michael Maier (1568-1622)/
Lang, Sarah. - Graz : Grazer Universitätsverlag - Leykam - Karl-Franzens-
Universität Graz, 2018\n    \n    \n    \n\n   \n\n    \n    2.\xa0\n
\n\nMichael Maier : nine newly discovered letters/ Lenke, Nils. - In: Ambix, Bd.
61 (2014), 1, S.1-47\n    \n    \n    \n\n   \n\n    \n    3.\xa0\n    \n\nDie
Tradierung alchemischen Wissens bei Michael Maier, Andreas Libavius und Oswald
Croll/ Wels, Volkhard. - In: Natur - Religion - Medien (2013), S.63-85\n    \n
\n    \n\n   \n\n    \n    4.\xa0\n    \n\nDoppelt verkettete Tricinien : Zarlino,
Calvisius und Michael Maier/ Braun, Werner. - In: Tempus musicæ - tempus mundi
(2008), S.103-116\n    \n    \n    \n\n   \n\n    \n    5.\xa0\n    \n\nMichael
Maiers Chymisches Cabinet : Atalanta fugiens deutsch nach der Ausgabe von 1708/
Maier, Michael. - Berlin [u.a.] : Thurneysser, 2007\n    \n    \n    \n\n   \n\n
\n    6.\xa0\n    \n\nOccult semiotics and iconology : Michael Maier\'s
alchemical emblems // Szönyi, György E.. - In: Mundus emblematicus (2003),
S.301-323\n    \n    \n    \n\n   \n\n    \n    7.\xa0\n    \n\nThe quest for the
Phoenix : spiritual alchemy and Rosicrucianism in the work of Count Michael
Maier (1569 - 1622)/ Tilton, Hereward. - Berlin [u.a.] : de Gruyter, 2003\n
\n    \n    \n\n   \n\n    \n    8.\xa0\n    \n\nArcana arcanissima : Emblematik
und Mythoalchemie bei Michael Maier/ Harzer, Friedmann. - In: Polyvalenz und
Multifunktionalität der Emblematik (2002), S.319-332\n    \n    \n    \n\n   \n\n
\n    9.\xa0\n    \n\nHermetische Poesie des Frühbarock : die "Cantilenae
intellectuales" Michael Maiers ; Edition mit Übersetzung, Kommentar und Bio-
Bibliographie/ Leibenguth, Erik. - Tübingen : Niemeyer, 2002\n    \n    \n    \n\n
\n\n    \n    10.\xa0\n    \n\nAtalante fugitive : Traduction française
d\'Étienne Perrot/ Maier, Michael. - Paris : Librairie de Médicis, 1970\n    \n
\n    \n\n    \n\n\n\n    \n    \n    \n\n    \n
\n\n1\xa0-\xa010\xa0von\xa069\xa0\xa0\xa0\xa0\xa0\xa0\n    \n    \n    \n    \n
\n\t\n\t\n\t\n\t\n    \n    \n    \n    \n    \xa0\n    \n\n    \n\n    \n
\n\n1\xa0-\xa010\xa0von\xa069\xa0\xa0\xa0\xa0\xa0\xa0\n    \n    \n    \n\n\t
\n    \n\n    \n\n\nvar coverLink = document.getElementById(\'coverLink\');\nvar

```
coverUrl  = document.getElementById(\'coverUrl\');\nif (coverUrl) coverUrl =
coverUrl.attributes.href.value;\nif (coverUrl && coverLink) {\n  var a =
document.createElement(\'a\');\n  a.setAttribute(\'href\',coverUrl);\n
a.setAttribute(\'target\',\'_blank\');\n  coverLink.appendChild(a);\n  var img =
document.createElement(\'img\');\n  img.addEventListener(\'load\', function() {
a.appendChild(img); });\n  if
(coverUrl.match(/(gbv\\.de\\/dms.+cov\\/[0-9X]+)\\.jpg$/)) {\n    coverUrl =
coverUrl.replace(/(cov\\/[0-9X]+)\\.jpg/,\'$1,400,400.jpg\');\n  }\n  img.src =
coverUrl.replace(/^https?:/,\'\');\n}\n\n\n'
```

## 2.1 The data is structured as follows

There is a table with `@summary` whose value is `hitlist`

```
table @summary="hitlist"
tbody
tr 1=picture
tr valign=top

td class=hit algin=left
```

–> text contains a link (the title) and the author name, first name - pub place: publisher, year.

```
a #InitialFocusPoint
href SHW?FRST=1
```

the results are then numbered from 1-69

## 2.2 Having fun

When we run the next cell, we realize that - oh my- my master's thesis shows up as the first search result. Well isn't this fun….

```
[7]: hits = tree.xpath("//table[@summary='hitlist']/tr[@valign='top']/
     ↪td[@class='hit' and @align='left']")
     for hit in hits:
         print_elem(hit)
```

```
<td>
Tags: {'class': 'hit', 'align': 'left', 'valign': 'top'};
  Ein religionswissenschaftlicher Kommentar zu den Arcana Arcanissima und der
Mythoalchemie des alchemo-hermetischen latrochemikers Michael Maier (1568-1622)/
Lang, Sarah. - Graz : Grazer Universitäts…


<td>
Tags: {'class': 'hit', 'align': 'left', 'valign': 'top'};
  Michael Maier : nine newly discovered letters/ Lenke, Nils. - In: Ambix, Bd.
```

61 (2014), 1, S.1-47      …


<td>
Tags: {'class': 'hit', 'align': 'left', 'valign': 'top'};
  Die Tradierung alchemischen Wissens bei Michael Maier, Andreas Libavius und
Oswald Croll/ Wels, Volkhard. - In: Natur - Religion - Medien (2013), S.63-85
…


<td>
Tags: {'class': 'hit', 'align': 'left', 'valign': 'top'};
  Doppelt verkettete Tricinien : Zarlino, Calvisius und Michael Maier/ Braun,
Werner. - In: Tempus musicæ - tempus mundi (2008), S.103-116      …


<td>
Tags: {'class': 'hit', 'align': 'left', 'valign': 'top'};
  Michael Maiers Chymisches Cabinet : Atalanta fugiens deutsch nach der Ausgabe
von 1708/ Maier, Michael. - Berlin [u.a.] : Thurneysser, 2007      …


<td>
Tags: {'class': 'hit', 'align': 'left', 'valign': 'top'};
  Occult semiotics and iconology : Michael Maier's alchemical emblems // Szönyi,
György E.. - In: Mundus emblematicus (2003), S.301-323      …


<td>
Tags: {'class': 'hit', 'align': 'left', 'valign': 'top'};
  The quest for the Phoenix : spiritual alchemy and Rosicrucianism in the work
of Count Michael Maier (1569 - 1622)/ Tilton, Hereward. - Berlin [u.a.] : de
Gruyter, 2003      …


<td>
Tags: {'class': 'hit', 'align': 'left', 'valign': 'top'};
  Arcana arcanissima : Emblematik und Mythoalchemie bei Michael Maier/ Harzer,
Friedmann. - In: Polyvalenz und Multifunktionalität der Emblematik (2002),
S.319-332      …


<td>
Tags: {'class': 'hit', 'align': 'left', 'valign': 'top'};
  Hermetische Poesie des Frühbarock : die "Cantilenae intellectuales" Michael
Maiers ; Edition mit Übersetzung, Kommentar und Bio-Bibliographie/ Leibenguth,
Erik. - Tübingen : Niemeyer, 2002      …

```
<td>
Tags: {'class': 'hit', 'align': 'left', 'valign': 'top'};
   Atalante fugitive : Traduction française d'Étienne Perrot/ Maier, Michael. -
Paris : Librairie de Médicis, 1970      …
```

## 3 Getting the data in structured form

Now that we have localized the data we wanted (yet didn't actually need..), we save it in structured from. * linklist has all the links (i.e. titles) * href attribute is the link to the catalog entry for the title, we want to save that just in case.

```python
titles = {}
linklist = tree.xpath("//table[@summary='hitlist']/tr[@valign='top']/
 ↪td[@class='hit' and @align='left']/a")
tds = tree.xpath("//table[@summary='hitlist']/tr[@valign='top']/td[@class='hit'␣
 ↪and @align='left']")

for link in linklist:
    print(link.text_content())
    print(link.attrib['href'])
```

```
Ein religionswissenschaftlicher Kommentar zu den Arcana Arcanissima und der
Mythoalchemie des alchemo-hermetischen latrochemikers Michael Maier (1568-1622)
SHW?FRST=1
Michael Maier : nine newly discovered letters
SHW?FRST=2
Die Tradierung alchemischen Wissens bei Michael Maier, Andreas Libavius und
Oswald Croll
SHW?FRST=3
Doppelt verkettete Tricinien : Zarlino, Calvisius und Michael Maier
SHW?FRST=4
Michael Maiers Chymisches Cabinet : Atalanta fugiens deutsch nach der Ausgabe
von 1708
SHW?FRST=5
Occult semiotics and iconology : Michael Maier's alchemical emblems /
SHW?FRST=6
The quest for the Phoenix : spiritual alchemy and Rosicrucianism in the work of
Count Michael Maier (1569 - 1622)
SHW?FRST=7
Arcana arcanissima : Emblematik und Mythoalchemie bei Michael Maier
SHW?FRST=8
Hermetische Poesie des Frühbarock : die "Cantilenae intellectuales" Michael
Maiers ; Edition mit Übersetzung, Kommentar und Bio-Bibliographie
```

SHW?FRST=9
Atalante fugitive : Traduction française d'Étienne Perrot
SHW?FRST=10

## 3.1   We can also get the whole citation like this:

```
[9]:  for td in tds:
          print(td.text_content())
```

Ein religionswissenschaftlicher Kommentar zu den Arcana Arcanissima und der
Mythoalchemie des alchemo-hermetischen Iatrochemikers Michael Maier (1568-1622)/
Lang, Sarah. – Graz : Grazer Universitätsverlag – Leykam – Karl-Franzens-
Universität Graz, 2018

Michael Maier : nine newly discovered letters/ Lenke, Nils. – In: Ambix, Bd. 61
(2014), 1, S.1-47

Die Tradierung alchemischen Wissens bei Michael Maier, Andreas Libavius und
Oswald Croll/ Wels, Volkhard. – In: Natur – Religion – Medien (2013), S.63-85

Doppelt verkettete Tricinien : Zarlino, Calvisius und Michael Maier/ Braun,
Werner. – In: Tempus musicæ – tempus mundi (2008), S.103-116

Michael Maiers Chymisches Cabinet : Atalanta fugiens deutsch nach der Ausgabe
von 1708/ Maier, Michael. – Berlin [u.a.] : Thurneysser, 2007

Occult semiotics and iconology : Michael Maier's alchemical emblems // Szönyi,
György E.. – In: Mundus emblematicus (2003), S.301-323

The quest for the Phoenix : spiritual alchemy and Rosicrucianism in the work of
Count Michael Maier (1569 – 1622)/ Tilton, Hereward. – Berlin [u.a.] : de
Gruyter, 2003

```
Arcana arcanissima : Emblematik und Mythoalchemie bei Michael Maier/ Harzer,
Friedmann. - In: Polyvalenz und Multifunktionalität der Emblematik (2002),
S.319-332
```

```
Hermetische Poesie des Frühbarock : die "Cantilenae intellectuales" Michael
Maiers ; Edition mit Übersetzung, Kommentar und Bio-Bibliographie/ Leibenguth,
Erik. - Tübingen : Niemeyer, 2002
```

```
Atalante fugitive : Traduction française d'Étienne Perrot/ Maier, Michael. -
Paris : Librairie de Médicis, 1970
```

- `http://opac.lbs-braunschweig.gbv.de/DB=2/SET=4/TTL=11/NXT?FRST=21` is the current page.
- The follow-up hrefs are SHW?FRST=num
- `http://opac.lbs-braunschweig.gbv.de/DB=2/SET=4/` + TTL=21/SHW?FRST=21

Using the link list generated earlier, we could do this for all the result pages so we'd get all the 69 results. And we could follow the follow-up links to the catalogue but the catalogue entries are not as easy to scrape, so we'll stop here.

[ ]: 

[ ]: