# Micro Archives as Rich Digital Object Representations

Helge Holzmann,
Mila Runnwerth
WebSci '18

# Referencing digital objects in scientific publications: Stable & dynamic data

| stable data | dynamic data |
|---|---|
| journal articles<br>set of measuring data | software<br>science blogs / forums |
|  |  |

# Referencing digital objects in scientific publications: The Web as infrastructure

Blog and software pages show the latest status of development or participation.

The Internet Archive unsystematically archives snapshots so web pages can be retraced over time.

Question: How complete and coherent is the archived web with respect to related resources linked on the corresponding web pages?

# Case Studies: Data sets

Blogs: TREC Blogs '08

- 28.488.766 permalinked blog posts
- from 1.011.733 homepages

Software: swMATH

- 21.947 software packages of mathematical software with corresponding URLs (May 29[th])

# Case Studies

**resources**
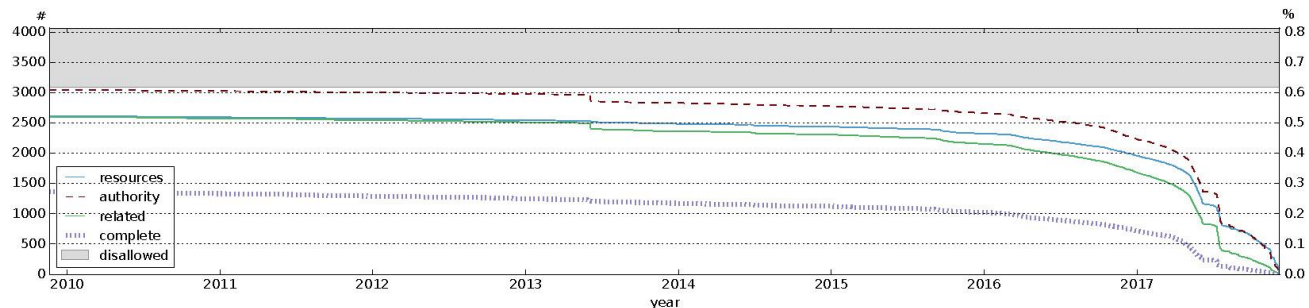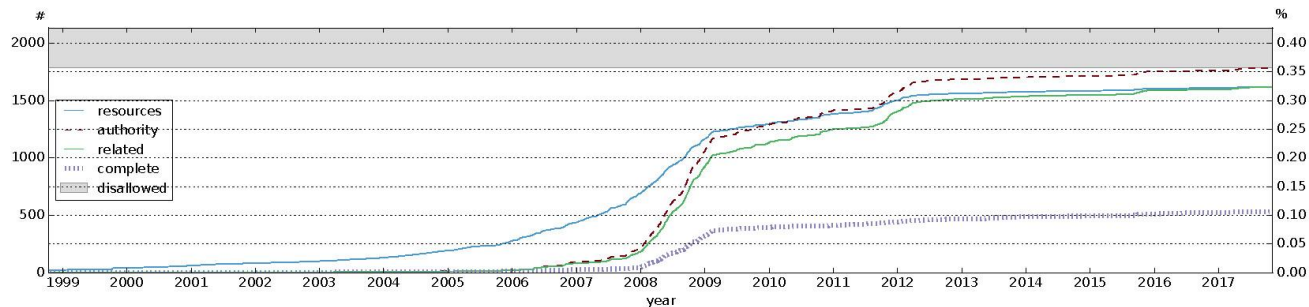an object as fraction of its archived resources

**authority**
authority pages only

**related**
fraction of resources only if *authority* is archived

**complete**

# Micro Archives Use Case Scenario

**Specifying Micro Archives**

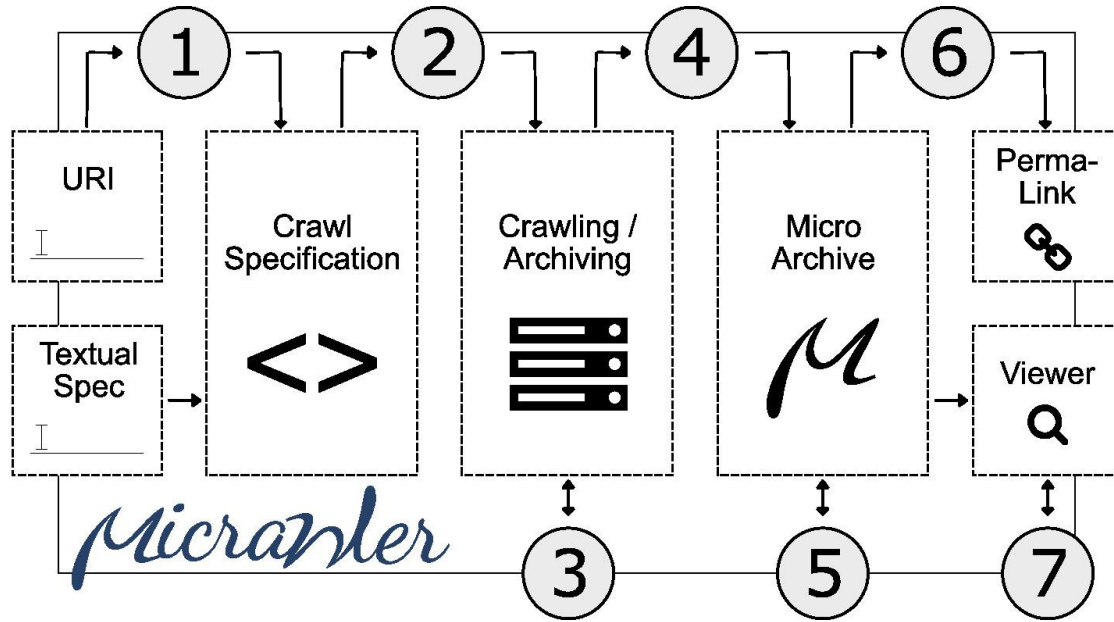Defining a digital representation of an object with certain specification

**Crawling / Archiving**

According to the specification crawling & archiving is executed at the same time in the specified depth, including the specified resources.

**Presentation / Citing**

The newly created *micro archive* is a semantic collection anchored in time that may be assigned a URL or, more specifically, a DOI for citing.

# Proof-of-concept prototype Micrawler



1. Spec Proxy
2. Crawl Queue
3. Archiving / Crawl Service
4. Archive Meta Service
5. Analysers
6. Persistence Provider
7. Viewer

# Proof-of-concept prototype Micrawler



Demo

# Tempas Micrawler

**What is a micro archive?** A micro archive is a snapshot of a fixed (evolving) set of URLs that are representative for some object or entity (at a given time). Hence, such an archive can be used to describe and / or derive information about its subject at the time of the crawl.

Create your own micro archive for an entity or object of your choice by either defining a set of URLs manually or loading / extracting a crawl specification from some URL:

**Enter spec definition**

Enter URL to load / extract a crawl specification

**Load / extract spec**

**Left window:**

Tempas Web Archive Search

*Tempas*
*Micrawler*

Please enter a crawl specification with one URL per line in the form: "Label: http://URL".

```
title=Helge Holzmann
type=person
# This micro crawl specification was extracted from http://www.helgeholzmann.de
Home:                     http://www.helgeholzmann.de
L3S Research Center:      http://www.L3S.de
Hannover, Germany:        https://www.google.de/maps/place/Hannover
My publications:          http://www.helgeholzmann.de/publications
helgeho on GitHub:        http://www.github.com/helgeho
@helgeho on Twitter:      https://twitter.com/search?q=%40helgeho
Contact:                  http://www.helgeholzmann.de/contact
Helge on arXiv:           http://www.arxiv.org/a/holzmann_h_1
Research:                 http://www.helgeholzmann.de/
```

[ Back ]  [ Show report ]  [ Start crawl ]

**Right window:**

Tempas Web Archive Search

*Tempas*
*Micrawler*

**Helge Holzmann** (person)
🕐 2018-05-11 07:55:24 UTC

show spec | re-crawl | cite

| **Home** | |
| http://www.helgeholzmann.de | 2018-05-11 07:55:26 UTC |

| **L3S Research Center** | |
| http://www.L3S.de | 2018-05-11 07:55:28 UTC |

| **Hannover, Germany** | |
| https://www.google.de/maps/place/Hannover | N/A |

| **My publications** | |
| http://www.helgeholzmann.de/publications | 2018-05-11 07:55:27 UTC |

| **helgeho on GitHub** | |
| http://www.github.com/helgeho | 2018-05-11 07:55:28 UTC |

| **@helgeho on Twitter** | |
| https://twitter.com/search?q=%40helgeho | 2018-05-11 07:55:26 UTC |

| **Contact** | |
| http://www.helgeholzmann.de/contact | 2018-05-01 17:33:57 UTC |

| **Helge on arXiv** | |
| http://www.arxiv.org/a/holzmann_h_1 | 2017-09-25 02:29:07 UTC |

Fork me on GitHub

# Tempas
## Micrawler
Temporal Web Archive Search

**What is a micro archive?** A micro archive is a snapshot of a fixed (evolving) set of URLs that are representative for some object or entity (at a given time). Hence, such an archive can be used to describe and / or derive information about its subject at the time of the crawl.

Create your own micro archive for an entity or object of your choice by either defining a set of URLs manually or loading / extracting a crawl specification from some URL:

Enter spec definition

Enter URL to load / extract a crawl specification

Load / extract spec

Tempas Web Archive Search

*Tempas*
*Micrawler*

**What is a micro archive?** A
that are representative for so
can be used to describe and
crawl.

Create your own micro archiv
set of URLs manually or loadi

Enter URL to load / extract

© 2017 L3S Resea

Tempas Web Archive Search

*Tempas*
*Micrawler*

Please enter a crawl specification with one URL per line in the form: "Label: http://URL".

```
Title=Dutch Olympic Swimming Team 2016

# My crawled homepages about the Dutch Olympic swimming team of 2016

General Wikipedia page about Swimming: https://en.wikipedia.org/wiki/Swimming_at_the_2016_Summ
General Wikipedia page about the Olympics 2016: https://en.wikipedia.org/wiki/Netherlands_at_t
Sebastiaan Verschuren:        https://en.wikipedia.org/wiki/Sebastiaan_Verschuren org,wikipedia,(
Dion Dreesens                 https://en.wikipedia.org/wiki/Dion_Dreesens org,wikipedia,en)/wil
Maarten Brzoskowski:          https://en.wikipedia.org/wiki/Maarten_Brzoskowski org,wikipedia,(
Joeri Verlinden:              https://en.wikipedia.org/wiki/Joeri_Verlinden org,wikipedia,en)/(
Ben Schwietert:               https://en.wikipedia.org/wiki/Ben_Schwietert org,wikipedia,en)/w
Kyle Stolk:                   https://nl.wikipedia.org/wiki/Kyle_Stolk org,wikipedia,nl)/wiki/l
Ferry Weertman:               https://en.wikipedia.org/wiki/Ferry_Weertman org,wikipedia,en)/wil
```

Back    Show report    Start crawl

© 2017 L3S Research Center (Helge Holzmann). All rights reserved.

Tempas Web Archive Search

**What is a micro archive?** A
that are representative for so
can be used to describe and
crawl.

Create your own micro archiv
set of URLs manually or loadi

Enter URL to load / extract

Tempas Web Archive Search

Please enter a crawl

Title=Dutch Olympic S

# My crawled homepage

General Wikipedia pag
General Wikipedia pag
Sebastiaan Verschuren
Dion Dreesens:
Maarten Brzoskowski:
Joeri Verlinden:
Ben Schwiertert:
Kyle Stolk:
Ferry Weertman:

Tempas Web Archive Search

**Please wait while we are capturing your micro crawl...**

12% - https://nl.wikipedia.org/wiki/Kyle_Stolk : done.
18% - https://en.wikipedia.org/wiki/Robin_Neumann : done.
25% - https://en.wikipedia.org/wiki/Femke_Heemskerk : done.
31% - https://en.wikipedia.org/wiki/Ferry_Weertman : done.
37% - https://en.wikipedia.org/wiki/Ranomi_Kromowidjojo : done.
43% - https://en.wikipedia.org/wiki/Ben_Schwiertert : done.
50% - https://en.wikipedia.org/wiki/Joeri_Verlinden : done.
56% - https://en.wikipedia.org/wiki/Dion_Dreesens : done.
62% - https://en.wikipedia.org/wiki/Maarten_Brzoskowski : done.
68% - https://en.wikipedia.org/wiki/Inge_Dekker : timeout.
75% - https://en.wikipedia.org/wiki/Sebastiaan_Verschuren : timeout.
81% - https://en.wikipedia.org/wiki/Netherlands_at_the_2016_Summer_Olympics : timeout.
87% - https://en.wikipedia.org/wiki/Swimming_at_the_2016_Summer_Olympics : timeout.
93% - Dutch : done.
100% - DutchOlympicSwimmingTeam2016 : done.
100% - crawl finished.

Back    Show report

# Tempas Micrawler

⊘ 2018-05-29 13:42:52 UTC

show spec | re-crawl | cite

| | |
|---|---|
| **General Wikipedia page about Swimming** | 2018-05-29 13:34:04 UTC |
| https://en.wikipedia.org/wiki/Swimming_at_the_2016_Summer_Olympics | |

| | |
|---|---|
| **General Wikipedia page about the Olympics 2016** | 2018-05-29 13:34:03 UTC |
| https://en.wikipedia.org/wiki/Netherlands_at_the_2016_Summer_Olympics | |

| | |
|---|---|
| **Sebastiaan Verschuren** | 2018-05-29 13:34:03 UTC |
| https://en.wikipedia.org/wiki/Sebastiaan_Verschuren | |

| | |
|---|---|
| **Dion Dreesens** | 2018-05-29 13:33:44 UTC |
| https://en.wikipedia.org/wiki/Dion_Dreesens | |

| | |
|---|---|
| **Maarten Brzoskowski** | 2018-05-29 13:33:44 UTC |
| https://en.wikipedia.org/wiki/Maarten_Brzoskowski | |

| | |
|---|---|
| **Joeri Verlinden** | 2018-05-29 13:33:44 UTC |
| https://en.wikipedia.org/wiki/Joeri_Verlinden | |

TemPoralas

Fork me on Git

Fork me on Git

## Cite Micro Crawl ✕

Permanent URL    BibTeX    BibLaTeX

```
@online{unnamed,
    url = {http://tempas.l3s.de/micrawler/permalink/785d426},
    urldate = {2018-05-29T13:50:53.000Z},
    note = {Archived using Micrawler}
}
```

Close

**General Wikipedia page about the Olympics 2016**    2018-05-29 13:51:18 UTC
https://en.wikipedia.org/wiki/Netherlands_at_the_2016_Summer_Olympics

**Sebastiaan Verschuren**    2018-05-29 13:51:10 UTC
https://en.wikipedia.org/wiki/Sebastiaan_Verschuren

**Dion Dreesens**    2018-05-29 13:50:57 UTC
https://en.wikipedia.org/wiki/Dion_Dreesens

**Maarten Brzoskowski**    2018-05-29 13:50:56 UTC
https://en.wikipedia.org/wiki/Maarten_Brzoskowski

**Joeri Verlinden**    2018-05-29 13:50:56 UTC
https://en.wikipedia.org/wiki/Joeri_Verlinden

**Joeri Verlinden**    2018-05-29 13:33:44 UTC
https://en.wikipedia.org/wiki/Joeri_Verlinden

# Outlook & opportunities

- Supporting web archives

- Temporally relevant collections

- Structuring the Web

- Rich information

holzmann@l3s.de          mila.runnwerth@tib.eu