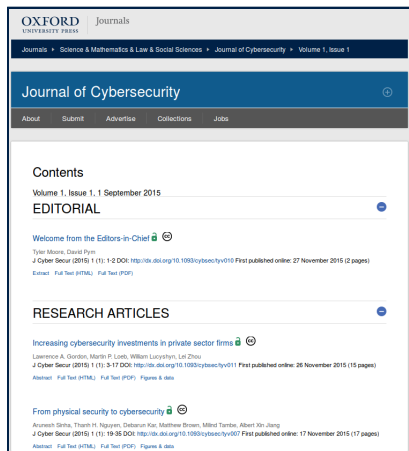# Smart Harvesting with OXPath

Christopher Michels
University of Trier
michelsc@uni-trier.de

Mandy Neumann
TH Köln
mandy.neumann@th-koeln.de

June 15, 2018

# Harvesting Bibliographic Data

# Accessing Bibliographic Data

# Accessing Bibliographic Data

# The Role of OXPath in Smart Harvesting II

Motivation:

- extract bibliographic data with OXPath

- facilitate maintenance of scientific literature databases

Solution:

- provide working environment and tools to use OXPath

OXPath:

- simple, declarative language for web data extraction

# Table of Contents

# Table of Contents

# Sources of Raw Web Data

# Sources of Raw Web Data

# Sources of Raw Web Data

# Sources of Raw Web Data

# Sources of Raw Web Data

# Sources of Raw Web Data

# Sources of Raw Web Data



Oxford Academic:

- OUP moved to new platform

- Winter 2016 - Spring 2017

- gradually moving individual journals

- 3 content platforms in use at the same time

# Sources of Raw Web Data

# Sources of Raw Web Data

# Sources of Raw Web Data

# Sources of Raw Web Data

# Sources of Raw Web Data

# Sources of Raw Web Data

# Sources of Raw Web Data

# Problem of Data Heterogeneity

Sources of raw bibliographic data vary largely in quality and format, e.g.:

- website layouts

- change

# Problem of Data Integration

- question of feasibility: automated vs. manual harvesting

- expensive maintenance

# DBLP as a Case Example

# DBLP as a Case Example

regular expressions

static-content handler

output formatter

DBLP API

Used across several steps, e.g.:

- publisher-key validation

- retrieving lists of issues

- retrieving tables of content

- retrieving records:
  `<tr[^>]*>.*?</tr>`

# DBLP as a Case Example

| DBLP API | OXPath API |
|---|---|
| regular expressions | OXPath expressions |
| static-content handler | controllable browser |
| output formatter | output-handler interface |

# Table of Contents

# What Is OXPath?

- simple, declarative language for web data extraction

- XPath extension:
    - actions
    - iteration
    - extraction

# What Is XPath?

- query language

- XML document as a tree of nodes

- XPath expressions as location paths

# What Is XPath?

```
C:\
│
├─ Program Files\
│   │
│   ├─ Atom
│   │
│   ├─ Eclipse
│   │
│   └─ Microsoft Office
│
└─ Users\
    │
    ├─ Jane Doe
    │
    └─ John Smith
```

### File-Path Examples

```
1 C:\Program Files\Microsoft Office
2 C:\Users\Jane Doe
```

# What Is XPath?

## Queried XML File

```xml
 1 <?xml version="1.0" encoding="UTF-8"?>
 2 <results>
 3   <record class="current">
 4     <volume>30</volume>
 5     <issue>11</issue>
 6     <year>2016</year>
 7     <url>http://.../tadr20/30/11</url>
 8   </record>
 9   <record>
10     <volume>30</volume>
11     <issue>10</issue>
12     <year>2016</year>
13     <url>http://.../tadr20/30/10</url>
14   </record>
15   <record>
16     <volume>30</volume>
17     <issue>9</issue>
18     <year>2016</year>
19     <url>http://.../tadr20/30/9</url>
20   </record>
21 </results>
```

## XPath Expression

```
1 /results/record/issue
```

## Result Set

# What Is XPath?

## Queried XML File

```xml
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <record class="current">
4     <volume>30</volume>
5     <issue>11</issue>
6     <year>2016</year>
7     <url>http://.../tadr20/30/11</url>
8   </record>
9   <record>
10    <volume>30</volume>
11    <issue>10</issue>
12    <year>2016</year>
13    <url>http://.../tadr20/30/10</url>
14  </record>
15  <record>
16    <volume>30</volume>
17    <issue>9</issue>
18    <year>2016</year>
19    <url>http://.../tadr20/30/9</url>
20  </record>
21 </results>
```

## XPath Expression

```
1 /results/record/issue
```

## Result Set

```
1 (
2   <issue>11</issue>,
3   <issue>10</issue>,
4   <issue>9</issue>
5 )
```

# What Is XPath?

## Queried XML File

```xml
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <record class="current">
4     <volume>30</volume>
5     <issue>11</issue>
6     <year>2016</year>
7     <url>http://.../tadr20/30/11</url>
8   </record>
9   <record>
10     <volume>30</volume>
11     <issue>10</issue>
12     <year>2016</year>
13     <url>http://.../tadr20/30/10</url>
14   </record>
15   <record>
16     <volume>30</volume>
17     <issue>9</issue>
18     <year>2016</year>
19     <url>http://.../tadr20/30/9</url>
20   </record>
21 </results>
```

## XPath Expression

```
1 /results/record/url/text()
```

## Result Set

# What Is XPath?

## Queried XML File

```xml
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <record class="current">
4     <volume>30</volume>
5     <issue>11</issue>
6     <year>2016</year>
7     <url>http://.../tadr20/30/11</url>
8   </record>
9   <record>
10    <volume>30</volume>
11    <issue>10</issue>
12    <year>2016</year>
13    <url>http://.../tadr20/30/10</url>
14  </record>
15  <record>
16    <volume>30</volume>
17    <issue>9</issue>
18    <year>2016</year>
19    <url>http://.../tadr20/30/9</url>
20  </record>
21 </results>
```

## XPath Expression

```
1 /results/record/url/text()
```

## Result Set

```
1 (
2   "http://.../toc/tadr20/30/11",
3   "http://.../toc/tadr20/30/10",
4   "http://.../toc/tadr20/30/9"
5 )
```

# What Is XPath?

## Queried XML File

```
 1 <?xml version="1.0" encoding="UTF-8"?>
 2 <results>
 3  <record class="current">
 4    <volume>30</volume>
 5    <issue>11</issue>
 6    <year>2016</year>
 7    <url>http://.../tadr20/30/11</url>
 8  </record>
 9  <record>
10    <volume>30</volume>
11    <issue>10</issue>
12    <year>2016</year>
13    <url>http://.../tadr20/30/10</url>
14  </record>
15  <record>
16    <volume>30</volume>
17    <issue>9</issue>
18    <year>2016</year>
19    <url>http://.../tadr20/30/9</url>
20  </record>
21 </results>
```

## XPath Expression

```
 1 /results/record[@class="current"]
```

## Result Set

# What Is XPath?

## Queried XML File

```
 1 <?xml version="1.0" encoding="UTF-8"?>
 2 <results>
 3   <record class="current">
 4     <volume>30</volume>
 5     <issue>11</issue>
 6     <year>2016</year>
 7     <url>http://.../tadr20/30/11</url>
 8   </record>
 9   <record>
10     <volume>30</volume>
11     <issue>10</issue>
12     <year>2016</year>
13     <url>http://.../tadr20/30/10</url>
14   </record>
15   <record>
16     <volume>30</volume>
17     <issue>9</issue>
18     <year>2016</year>
19     <url>http://.../tadr20/30/9</url>
20   </record>
21 </results>
```

## XPath Expression

```
 1 /results/record[@class="current"]
```

## Result Set

```
 1 (
 2   <record class="current">
 3     <volume>30</volume>
 4     <issue>11</issue>
 5     <year>2016</year>
 6     <url>[...]</url>
 7   </record>
 8 )
```

# What Does OXPath Add?

Action:

- fill in forms

- click links, buttons, etc.

Extraction:

- add markers to extract selected nodes

Iteration:

- loops, e.g. for paginated content

# Example: Navigating Google Scholar



## OXPath Expression

```
1 doc('https://scholar.google.com')
```

# Example: Navigating Google Scholar



**OXPath Expression**

```
1 doc('https://scholar.google.com')
2   //input[@id='gs_hdr_tsi']/{"OXPath"}
```

# Example: Navigating Google Scholar



## OXPath Expression

```
1 doc('https://scholar.google.com')
2 //input[@id='gs_hdr_tsi']/{"OXPath"}
3 /../following-sibling::button/{click/}
```

# Example: Navigating Google Scholar



## OXPath Expression

```
1 doc('https://scholar.google.com')
2 //input[@id='gs_hdr_tsi']/{"OXPath"}
3 ../following-sibling::button/{click/}
```

# Example: Navigating Google Scholar



## OXPath Expression

```
1 doc('https://scholar.google.com')
2   //input[@id='gs_hdr_tsi']/{"OXPath"}
3   /../following-sibling::button/{click/}
4     //*[@id='gs_res_ab_yy-b']/{click}
```

# Example: Navigating Google Scholar



OXPath Expression

```
1 doc('https://scholar.google.com')
2   //input[@id='gs_hdr_tsi']/{"OXPath"}
3   /../following-sibling::button/{click/}
4   //*[@id='gs_res_ab_yy-b']/{click/}
5     //following::*[@role='menuitemradio'][contains(.,
         '2016')]/{click/}
```

# Example: Navigating Google Scholar



## OXPath Expression

```
1 doc('https://scholar.google.com')
2  //input[@id='gs_hdr_tsi']/{"OXPath"}
3  /../following-sibling::button/{click/}
4   //*[@id='gs_res_ab_yy-b']/{click/}
5    //following::*[@role='menuitemradio'][contains(.,
        '2016')]/{click/}
```

# Example: Navigating Google Scholar



## OXPath Expression

```
1 doc('https://scholar.google.com')
2   //input[@id='gs_hdr_tsi']/{"OXPath"}
3   /../following-sibling::button/{click/}
4     //*[@id='gs_res_ab_yy-b']/{click}
5       //following::*[@role='menuitemradio'][contains(.,
            '2016')]/{click/}
6   //div[@class='gs_ri']//h3/a:<title=string(.)>
```

## XML Output

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <title>Tim Furche, Georg Gottlob, [...]</title>
4 </results>
```

# Example: Navigating Google Scholar



## OXPath Expression

```
1 doc('https://scholar.google.com')
2 //input[@id='gs_hdr_tsi']/{"OXPath"}
3 /../following-sibling::button/{click/}
4   //*[@id='gs_res_ab_yy-b']/{click}
5     //following::*[@role='menuitemradio'][contains(.,
        '2016')]/{click/}
6 /(//*[@id='gs_nm']/button[2][not(@disabled)]/{click/})*
7   //div[@class='gs_ri']//h3/a:<title=string(.)>
```

## XML Output

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <title>Tim Furche, Georg Gottlob, [...]</title>
4 </results>
```

# Example: Navigating Google Scholar



## OXPath Expression

```
1 doc('https://scholar.google.com')
2 //input[@id='gs_hdr_tsi']/{"OXPath"}
3 /../following-sibling::button/{click/}
4   //*[@id='gs_res_ab_yy-b']/{click}
5     //following::*[@role='menuitemradio'][contains(.,
      '2016')]/{click/}
6 /(//*[@id='gs_nm']/button[2][not(@disabled)]/{click/})*
7   //div[@class='gs_ri']//h3/a:<title=string(.)>
```

## XML Output

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <title>Tim Furche, Georg Gottlob, [...]</title>
4 </results>
```

# Example: Navigating Google Scholar



## OXPath Expression

```
1 doc('https://scholar.google.com')
2 //input[@id='gs_hdr_tsi']/{"OXPath"}
3 /../following-sibling::button/{click/}
4   //*[@id='gs_res_ab_yy-b']/{click}
5   //following::*[@role='menuitemradio'][contains(.,
      '2016')]/{click/}
6 /(//*[@id='gs_nm']/button[2][not(@disabled)]/{click/})*
7   //div[@class='gs_ri']//h3/a:<title=string(.)>
```

## XML Output

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <title>Tim Furche, Georg Gottlob, [...]</title>
4   <title>Special Issue: Big Data [...]</title>
5 </results>
```

# Example: Navigating Google Scholar



**OXPath Expression**

```
1 doc('https://scholar.google.com')
2 //input[@id='gs_hdr_tsi']/{"OXPath"}
3 /../following-sibling::button/{click/}
4   //*[@id='gs_res_ab_yy-b']/{click}
5     //following::*[@role='menuitemradio'][contains(.,
        '2016')]/{click/}
6 /(//*[@id='gs_nm']/button[2][not(@disabled)]/{click/})*
7   //div[@class='gs_ri']//h3/a:<title=string(.)>
```

**XML Output**

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <title>Tim Furche, Georg Gottlob, [...]</title>
4   <title>Special Issue: Big Data [...]</title>
5   <!--[...]-->
6 </results>
```

# Example: Navigating Google Scholar



## OXPath Expression

```
1 doc('https://scholar.google.com')
2   //input[@id='gs_hdr_tsi']/{"OXPath"}
3   /../following-sibling::button/{click/}
4     //*[@id='gs_res_ab_yy-b']/{click}
5       //following::*[@role='menuitemradio'][contains(.,
          '2016')]/{click/}
6   /(//*[@id='gs_nm']/button[2][not(@disabled)]/{click/})*
7     //div[@class='gs_ri']//h3/a:<title=string(.)>
```

## XML Output

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <title>Tim Furche, Georg Gottlob, [...]</title>
4   <title>Special Issue: Big Data [...]</title>
5   <!--[...]-->
6 </results>
```

# Why OXPath?

| XPath | OXPath |
|---|---|
| | |

- static web
- plain HTML
- complete content

- dynamic web
- AJAX
- content on demand

# OXPath In Use

In current working environment and production:

- OXPath-based wrappers in dblp

- collection of OXPath expressions prototypical of bibliographic domain

- OXPath Tutorial

# OXPath Tutorial



**Introduction to OXPath**

Ruslan R. Fayzrakhmanov
Christopher Michels
Mandy Neumann

- supported by the Oxford research team that developed OXPath

- including:
  - a concise summary of XPath
  - setup and use instructions for OXPath
  - list of action keywords
  - list of functions for extraction and DOM-tree navigation
  - starter examples from the bibliographic domain

# OXPath In Use

In current working environment and production:

- OXPath-based wrappers in dblp

- collection of OXPath expressions prototypical of bibliographic
  domain

- OXPath Tutorial

- integrated support for OXPath language in open-source editor *Atom*

# Tool Support: Atom

Language plugin for Atom text editor

- syntax highlighting for keywords

- helps spotting errors and improves readability

- intended to lower barriers for beginners

```
1  doc('https://scholar.google.com')
2  //input[@id='gs_hdr_tsi']/{"OXPath"}
3  /../following-sibling::button/{click/}
4  //*[@id='gs_res_ab_yy-b']/{click}
5  //following::*[@role='menuitemradio'][contains(.,'2016')]/{click/}
6  /(//*[@id='gs_nm']/button[2][not(@disabled)]/{click/})*
7  //div[@class='gs_ri']//h3/a:<title=string(.)>
8
```

```
1  doc('https://scholar.google.com')
2  //input[@id='gs_hdr_tsi']/{"OXPath"}
3  /../following-sibling::button/{click/}
4  //*[@id='gs_res_ab_yy-b']/{click}
5  //following::*[@role='menuitemradio'][contains(.,'2016')]/{click/}
6  /(//*[@id='gs_nm']/button[2][not(@disabled)]/{click/})*
7  //div[@class='gs_ri']//h3/a:<title=string(.)>
8
```

# OXPath In Use

In current working environment and production:

- OXPath-based wrappers in dblp

- collection of OXPath expressions prototypical of bibliographic domain

- OXPath Tutorial

- integrated support for OXPath language in open-source editor *Atom*

- OXPath-based monitoring for digital libraries of large-scale publishers

# Monitoring: ACM



Recently loaded issues and proceedings:
*(available in the DL within the past 2 weeks)*

Proceedings of the 10th International Conference on Security of Information and Networks
SIN '17

Proceedings of the 12th International Workshop on Variability Modelling of Software-Intensive Systems
VAMOS 2018

Proceedings of the 15th International Conference on Advances in Mobile Computing & Multimedia
MoMM2017

Proceedings of the 1st Reversing and Offensive-oriented Trends Symposium

- metadata delivery might be unreliable (e.g. incomplete)

- observe several news windows, unreliable as well
  - last 2 weeks

# Monitoring: ACM



- metadata delivery might be unreliable (e.g. incomplete)

- observe several news windows, unreliable as well
  - last 2 weeks
  - last 12 months

# Monitoring: ACM



- metadata delivery might be unreliable (e.g. incomplete)

- observe several news windows, unreliable as well
    - last 2 weeks

    - last 12 months

    - last 3 months

# OXPath In Use

In future use:

- interactive tool for devising OXPath expressions

- integrate further tools for OXPath

# Visual OXPath

- provide a visual interface
- integrate tools for XPath
- goal: semi-supervised wrapper induction

# Discussion

Thank you for your attention!
Feel free to ask any questions now!

Contact us:
`mandy.neumann@th-koeln.de`
`michelsc@uni-trier.de`

Source:
Visit `http://www.oxpath.org`, e.g.
to find the tutorial

# Table of Contents

# EDM 2014: Simple Extraction



## OXPath Expression

```
1 doc('http://edm2014.org/?page=proceedings')
```

## XML Output

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
1 </results>
```

# EDM 2014: Simple Extraction

## HTML Source

```html
1 <html xmlns="[...]" xml:lang="en">
2   <!--[...]-->
3   <div id="content">
4     <!--[...]-->
5     <strong>Online Proceedings</strong>
6     <!--[...]-->
7     <strong>Full Papers</strong>
8     <!--[...]-->
9     <p>Adaptive Practice of [...]
10       <br/>
11       <em> Jan Papousek, [...]</em>
12       <br/>
13       Pages 6-13 [
14       <a href="uploads/[...].pdf">pdf</a>
15       ]
16       <!--[...]-->
17     </p>
18     <!--[...]-->
19   </div>
20   <!--[...]-->
21 </html>
```

## OXPath Expression

```
1 doc('http://edm2014.org/?page=proceedings')
```

## XML Output

```xml
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
1 </results>
```

# EDM 2014: Simple Extraction

## HTML Source

```
1 <html xmlns="[...]" xml:lang="en">
2    <!--[...]-->
3    <div id="content">
4       <!--[...]-->
5       <strong>Online Proceedings</strong>
6       <!--[...]-->
7       <strong>Full Papers</strong>
8       <!--[...]-->
9       <p>Adaptive Practice of [...]
10         <br/>
11         <em> Jan Papousek, [...]</em>
12         <br/>
13         Pages 6-13 [
14         <a href="uploads/[...].pdf">pdf</a>
15         ]
16         <!--[...]-->
17      </p>
18      <!--[...]-->
19   </div>
20   <!--[...]-->
21 </html>
```

## OXPath Expression

```
1 doc('http://edm2014.org/?page=proceedings')
2    //*[@id='content']/p[./em]:<record>
```

## XML Output

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3    <record></record>
4    <record></record>
5    <!--[...]-->
6 </results>
```

# EDM 2014: Simple Extraction

## HTML Source

```
 1 <html xmlns="[...]" xml:lang="en">
 2   <!--[...]-->
 3   <div id="content">
 4     <!--[...]-->
 5     <strong>Online Proceedings</strong>
 6     <!--[...]-->
 7     <strong>Full Papers</strong>
 8     <!--[...]-->
 9     <p>Adaptive Practice of [...]
10       <br/>
11       <em> Jan Papousek, [...]</em>
12       <br/>
13       Pages 6-13 [
14       <a href="uploads/[...].pdf">pdf</a>
15       ]
16       <!--[...]-->
17     </p>
18     <!--[...]-->
19   </div>
20   <!--[...]-->
21 </html>
```

## OXPath Expression

```
 1 doc('http://edm2014.org/?page=proceedings')
 2 //*[@id='content']/p[./em]:<record>
 3   [./em:<authors=string(.)>]
```

## XML Output

```
 1 <?xml version="1.0" encoding="UTF-8"?>
 2 <results>
 3   <!--[...]-->
 4   <record>
 5     <authors> Jan Papousek, [...]</authors>
 6   </record>
 7   <!--[...]-->
 8 </results>
```

# EDM 2014: Simple Extraction

## HTML Source

```
 1 <html xmlns="[...]" xml:lang="en">
 2   <!--[...]-->
 3   <div id="content">
 4     <!--[...]-->
 5     <strong>Online Proceedings</strong>
 6     <!--[...]-->
 7     <strong>Full Papers</strong>
 8     <!--[...]-->
 9     <p>Adaptive Practice of [...]
10       <br/>
11       <em> Jan Papousek, [...]</em>
12       <br/>
13       Pages 6-13 [
14       <a href="uploads/[...].pdf">pdf</a>
15       ]
16       <!--[...]-->
17     </p>
18     <!--[...]-->
19   </div>
20   <!--[...]-->
21 </html>
```

## OXPath Expression

```
 1 doc('http://edm2014.org/?page=proceedings')
 2 //*[@id='content']/p[./em]:<record>
 3   [./em:<authors=string(.)>]
 4   [./text()[1]:<title=string(.)>]
```

## XML Output

```
 1 <?xml version="1.0" encoding="UTF-8"?>
 2 <results>
 3   <!--[...]-->
 4   <record>
 5     <authors> Jan Papousek, [...]</authors>
 6     <title>Adaptive Practice of [...]</title>
 7   </record>
 8   <!--[...]-->
 9 </results>
```

# EDM 2014: Simple Extraction

## HTML Source

```
 1 <html xmlns="[...]" xml:lang="en">
 2   <!--[...]-->
 3   <div id="content">
 4     <!--[...]-->
 5     <strong>Online Proceedings</strong>
 6     <!--[...]-->
 7     <strong>Full Papers</strong>
 8     <!--[...]-->
 9     <p>Adaptive Practice of [...]
10       <br/>
11       <em> Jan Papousek, [...]</em>
12       <br/>
13       Pages 6-13 [
14       <a href="uploads/[...].pdf">pdf</a>
15       ]
16     <!--[...]-->
17   </p>
18   <!--[...]-->
19   </div>
20   <!--[...]-->
21 </html>
```

## OXPath Expression

```
1 doc('http://edm2014.org/?page=proceedings')
2 //*[@id='content']/p[./em]:<record>
3   [./em:<authors=string(.)>]
4   [./text()[1]:<title=string(.)>]
5   [./br[2]/following-sibling::text()[1]
6       :<pages=substring-after(., "Pages ")>]
```

## XML Output

```
 1 <?xml version="1.0" encoding="UTF-8"?>
 2 <results>
 3   <!--[...]-->
 4   <record>
 5     <authors> Jan Papousek, [...]</authors>
 6     <title>Adaptive Practice of [...]</title>
 7     <pages>6-13 [</pages>
 8   </record>
 9   <!--[...]-->
10 </results>
```

# EDM 2014: Simple Extraction

## HTML Source

```
1 <html xmlns="[...]" xml:lang="en">
2   <!--[...]-->
3   <div id="content">
4     <!--[...]-->
5     <strong>Online Proceedings</strong>
6     <!--[...]-->
7     <strong>Full Papers</strong>
8     <!--[...]-->
9     <p>Adaptive Practice of [...]
10      <br/>
11      <em> Jan Papousek, [...]</em>
12      <br/>
13      Pages 6-13 [
14      <a href="uploads/[...].pdf">pdf</a>
15      ]
16      <!--[...]-->
17    </p>
18    <!--[...]-->
19  </div>
20  <!--[...]-->
21 </html>
```

## OXPath Expression

```
1 doc('http://edm2014.org/?page=proceedings')
2 //*[@id='content']/p[./em]:<record>
3   [./em:<authors=string(.)>]
4   [./text()[1]:<title=string(.)>]
5   [./br[2]/following-sibling::text()[1]
6      :<pages=substring-after(., "Pages ")>]
7   [./a:<url=string(@href)>]
```

## XML Output

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <!--[...]-->
4   <record>
5     <authors> Jan Papousek, [...]</authors>
6     <title>Adaptive Practice of [...]</title>
7     <pages>6-13 [</pages>
8     <url>uploads/[...].pdf</url>
9   </record>
10  <!--[...]-->
11 </results>
```

# EDM 2014: Simple Extraction

## HTML Source

```
1 <html xmlns="[...]" xml:lang="en">
2   <!--[...]-->
3   <div id="content">
4     <!--[...]-->
5     <strong>Online Proceedings</strong>
6     <!--[...]-->
7     <strong>Full Papers</strong>
8     <!--[...]-->
9     <p>Adaptive Practice of [...]
10      <br/>
11      <em> Jan Papousek, [...]</em>
12      <br/>
13      Pages 6-13 [
14      <a href="uploads/[...].pdf">pdf</a>
15      ]
16      <!--[...]-->
17    </p>
18    <!--[...]-->
19  </div>
20  <!--[...]-->
21 </html>
```

## OXPath Expression

```
1 doc('http://edm2014.org/?page=proceedings')
2 //*[@id='content']/p[./em]:<record>
3   [./em:<authors=string(.)>]
4   [./text()[1]:<title=string(.)>]
5   [./br[2]/following-sibling::text()[1]
6        :<pages=substring-after(., "Pages ")>]
7   [./a:<url=string(@href)>]
8   [./preceding::strong[1]:<header=string(.)>]
```

## XML Output

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <!--[...]-->
4   <record>
5     <authors> Jan Papousek, [...]</authors>
6     <title>Adaptive Practice of [...]</title>
7     <pages>6-13 [</pages>
8     <url>uploads/[...].pdf</url>
9     <header>Full Papers</header>
10  </record>
11  <!--[...]-->
12 </results>
```

# EDM 2014: Advanced Extraction

## HTML Source

```
1 <html xmlns="[...]" xml:lang="en">
2   <!--[...]-->
3   <div id="content">
4     <!--[...]-->
5     <strong>Online Proceedings</strong>
6     <!--[...]-->
7     <strong>Full Papers</strong>
8     <!--[...]-->
9     <p>Adaptive Practice of [...]
10       <br/>
11       <em> Jan Papousek, [...]</em>
12       <br/>
13       Pages 6-13 [
14       <a href="uploads/[...].pdf">pdf</a>
15       ]
16     <!--[...]-->
17     </p>
18     <!--[...]-->
19   </div>
20   <!--[...]-->
21 </html>
```

## OXPath Expression

```
1 doc('http://edm2014.org/?page=proceedings')
2 //*[@id='content']/p[./em]:<record>
3   [./em:<authors=string(.)>]
4   [./text()[1]:<title=string(.)>]
5   [./br[2]/following-sibling::text()[1]
       :<pages=substring-after(., "Pages ")>]
6   [./a:<url=string(@href)>]
7   [./preceding::strong[1]:<header=string(.)>]
```

## XML Output

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <!--[...]-->
4   <record>
5     <authors> Jan Papousek, [...]</authors>
6     <title>Adaptive Practice of [...]</title>
7     <pages>6-13 [</pages>
8     <url>uploads/[...].pdf</url>
9     <header>Full Papers</header>
10  </record>
11  <!--[...]-->
12 </results>
```

# EDM 2014: Advanced Extraction

## HTML Source

```
1 <html xmlns="[...]" xml:lang="en">
2   <!--[...]-->
3   <div id="content">
4     <!--[...]-->
5     <strong>Online Proceedings</strong>
6     <!--[...]-->
7     <strong>Full Papers</strong>
8     <!--[...]-->
9     <p>Adaptive Practice of [...]
10      <br/>
11      <em> Jan Papousek, [...]</em>
12      <br/>
13      Pages 6-13 [
14      <a href="uploads/[...].pdf">pdf</a>
15      ]
16      <!--[...]-->
17    </p>
18    <!--[...]-->
19  </div>
20  <!--[...]-->
21 </html>
```

## OXPath Expression

```
1 doc('http://edm2014.org/?page=proceedings')
2 //*[@id='content']/p[./em]:<record>
3   [./em:<authors=normalize-space(.)>]
4   [./text()[1]:<title=string(.)>]
5   [./br[2]/following-sibling::text()[1]
         :<pages=substring-after(., "Pages ")>]
6   [./a:<url=string(@href)>]
7   [./preceding::strong[1]:<header=string(.)>]
```

## XML Output

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <!--[...]-->
4   <record>
5     <authors>Jan Papousek, [...]</authors>
6     <title>Adaptive Practice of [...]</title>
7     <pages>6-13 [</pages>
8     <url>uploads/[...].pdf</url>
9     <header>Full Papers</header>
10  </record>
11  <!--[...]-->
12 </results>
```

# EDM 2014: Advanced Extraction

## HTML Source

```
1 <html xmlns="[...]" xml:lang="en">
2   <!--[...]-->
3   <div id="content">
4     <!--[...]-->
5     <strong>Online Proceedings</strong>
6     <!--[...]-->
7     <strong>Full Papers</strong>
8     <!--[...]-->
9     <p>Adaptive Practice of [...]
10      <br/>
11      <em> Jan Papousek, [...]</em>
12      <br/>
13      Pages 6-13 [
14      <a href="uploads/[...].pdf">pdf</a>
15      ]
16      <!--[...]-->
17    </p>
18    <!--[...]-->
19  </div>
20  <!--[...]-->
21 </html>
```

## OXPath Expression

```
1 doc('http://edm2014.org/?page=proceedings')
2 //*[@id='content']/p[./em]:<record>
3   [./em:<authors=normalize-space(.)>]
4   [./text()[1]:<title=string(.)>]
5   [./br[2]/following-sibling::text()[1]
       :<pages=replace(normalize-space(.),
       ".*?(\d+(-\d+)?).*", "$1")>]
6   [./a:<url=string(@href)>]
7   [./preceding::strong[1]:<header=string(.)>]
```

## XML Output

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <!--[...]-->
4   <record>
5     <authors>Jan Papousek, [...]</authors>
6     <title>Adaptive Practice of [...]</title>
7     <pages>6-13</pages>
8     <url>uploads/[...].pdf</url>
9     <header>Full Papers</header>
10  </record>
11  <!--[...]-->
12 </results>
```

# EDM 2014: Advanced Extraction

## HTML Source

```html
1 <html xmlns="[...]" xml:lang="en">
2   <!--[...]-->
3   <div id="content">
4     <!--[...]-->
5     <strong>Online Proceedings</strong>
6     <!--[...]-->
7     <strong>Full Papers</strong>
8     <!--[...]-->
9     <p>Adaptive Practice of [...]
10      <br/>
11      <em> Jan Papousek, [...]</em>
12      <br/>
13      Pages 6-13 [
14      <a href="uploads/[...].pdf">pdf</a>
15      ]
16      <!--[...]-->
17    </p>
18    <!--[...]-->
19  </div>
20  <!--[...]-->
21 </html>
```

## OXPath Expression

```
1 doc('http://edm2014.org/?page=proceedings')
2 //*[@id='content']/p[./em]:<record>
3   [./em:<authors=normalize-space(.)>]
4   [./text()[1]:<title=string(.)>]
5   [./br[2]/following-sibling::text()[1]
6      :<pages=replace(normalize-space(.),
7      ".*?(\d+(-\d+)?).*", "$1")>]
8   [./a:<url=qualify-url(@href)>]
9   [./preceding::strong[1]:<header=string(.)>]
```

## XML Output

```xml
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <!--[...]-->
4   <record>
5     <authors>Jan Papousek, [...]</authors>
6     <title>Adaptive Practice of [...]</title>
7     <pages>6-13</pages>
8     <url>http://[...]uploads/[...].pdf</url>
9     <header>Full Papers</header>
10  </record>
11  <!--[...]-->
12 </results>
```

# EDM 2016: Adapting from EDM 2014

## HTML Source (EDM 2014)

```
 1 <html xmlns="[...]" xml:lang="en">
 2   <!--[...]-->
 3   <div id="content">
 4     <!--[...]-->
 5     <strong>Online Proceedings</strong>
 6     <!--[...]-->
 7     <strong>Full Papers</strong>
 8     <!--[...]-->
 9     <p>Adaptive Practice of [...]
10       <br/>
11       <em> Jan Papousek, [...]</em>
12       <br/>
13       Pages 6-13 [
14       <a href="uploads/[...].pdf">pdf</a>
15       ]
16       <!--[...]-->
17     </p>
18     <!--[...]-->
19   </div>
20   <!--[...]-->
21 </html>
```

## OXPath Expression (EDM 2014)

```
 1 doc('http://edm2014.org/?page=proceedings')
 2 //*[@id='content']/p[./em]:<record>
 3   [./em:<authors=string(.)>]
 4   [./text()[1]:<title=string(.)>]
 5   [./br[2]/following-sibling::text()[1]
 6     :<pages=substring-after(., "Pages ")>]
 7   [./a:<url=string(@href)>]
 8   [./preceding::strong[1]:<header=string(.)>]
```

## XML Output (EDM 2014)

```
 1 <?xml version="1.0" encoding="UTF-8"?>
 2 <results>
 3   <!--[...]-->
 4   <record>
 5     <authors> Jan Papousek, [...]</authors>
 6     <title>Adaptive Practice of [...]</title>
 7     <pages>6-13 [</pages>
 8     <url>uploads/[...].pdf</url>
 9     <header>Full Papers</header>
10   </record>
11   <!--[...]-->
12 </results>
```

# EDM 2016: Adapting from EDM 2014



## OXPath Expression

```
1 doc('http://edm2016.org/proceedings.html')
2 //*[@id='content']/p[./em]:<record>
3    [./em:<authors=string(.)>]
4    [./text()[1]:<title=string(.)>]
5    [./br[2]/following-sibling::text()[1]
        :<pages=substring-after(., "Pages ")>]
6    [./a:<url=string(@href)>]
7    [./preceding::strong[1]:<header=string(.)>]
```

## XML Output

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3 </results>
```

# EDM 2016: Adapting from EDM 2014



## OXPath Expression

```
1  doc('http://edm2016.org/proceedings.html')
2  //*[@id='content']/p[./em]:<record>
3    [./em:<authors=string(.)>]
4    [./text()[1]:<title=string(.)>]
5    [./br[2]/following-sibling::text()[1]
       :<pages=substring-after(., "Pages ")>]
6    [./a:<url=string(@href)>]
7    [./preceding::strong[1]:<header=string(.)>]
```

## XML Output

```
1  <?xml version="1.0" encoding="UTF-8"?>
2  <results>
3    <!--[...]-->
4    <record>
5      <authors>???</authors>
6      <title>???</title>
7      <pages>???</pages>
8      <url>???</url>
9      <header>???</header>
10   </record>
11   <!--[...]-->
12 </results>
```

# EDM 2016: Adapting from EDM 2014



## OXPath Expression

```
1 doc('http://edm2016.org/proceedings.html')
2 ???:<record>
3   [???:<authors=???>]
4   [???:<title=???>]
5   [???:<pages=???>]
6   [???:<url=???>]
7   [???:<header=???>]
```

## XML Output

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <!--[...]-->
4   <record>
5     <authors>???</authors>
6     <title>???</title>
7     <pages>???</pages>
8     <url>???</url>
9     <header>???</header>
10  </record>
11  <!--[...]-->
12 </results>
```

# EDM 2016: Adapting from EDM 2014

## HTML Source

```
 1 <html>
 2   <!--[...]-->
 3   <h1>Individual papers</h1>
 4   <h3>Invited Talks</h3>
 5   <p>
 6     <a id="[...]" class="citation_title"
          href="[...]">Data-Driven [...]</a>
 7     <!--[...]-->
 8     <span class="[...]title">9th [...]</span>
 9     <span class="[...]firstpage">2</span>
10     <span class="[...]lastpage">2</span>
11     <span class="[...]pdf_url">http[...]</span>
12     <br/>
13     <span class="[...]author">Ra[...]</span>
14     <!--[...]-->
15   </p>
16   <!--[...]-->
17 </html>
```

## OXPath Expression

```
1 doc('http://edm2016.org/proceedings.html')
2 ???:<record>
3   [???:<authors=???>]
4   [???:<title=???>]
5   [???:<pages=???>]
6   [???:<url=???>]
7   [???:<header=???>]
```

## XML Output

```
 1 <?xml version="1.0" encoding="UTF-8"?>
 2 <results>
 3   <!--[...]-->
 4   <record>
 5     <authors>???</authors>
 6     <title>???</title>
 7     <pages>???</pages>
 8     <url>???</url>
 9     <header>???</header>
10   </record>
11   <!--[...]-->
12 </results>
```

# EDM 2016: Adapting from EDM 2014

## HTML Source

```
1  <html>
2    <!--[...]-->
3    <h1>Individual papers</h1>
4    <h3>Invited Talks</h3>
5    <p>
6      <a id="[...]" class="citation_title"
            href="[...]">Data-Driven [...]</a>
7      <!--[...]-->
8      <span class="[...]title">9th [...]</span>
9      <span class="[...]firstpage">2</span>
10     <span class="[...]lastpage">2</span>
11     <span class="[...]pdf_url">http[...]</span>
12     <br/>
13     <span class="[...]author">Ra[...]</span>
14     <!--[...]-->
15   </p>
16   <!--[...]-->
17 </html>
```

## OXPath Expression

```
1 doc('http://edm2016.org/proceedings.html')
2 //p[./*[contains(@class, 'cit')]]:<record>
3   [???:<authors=???>]
4   [???:<title=???>]
5   [???:<pages=???>]
6   [???:<url=???>]
7   [???:<header=???>]
```

## XML Output

```
1  <?xml version="1.0" encoding="UTF-8"?>
2  <results>
3    <!--[...]-->
4    <record>
5      <authors>???</authors>
6      <title>???</title>
7      <pages>???</pages>
8      <url>???</url>
9      <header>???</header>
10   </record>
11   <!--[...]-->
12 </results>
```

# EDM 2016: Adapting from EDM 2014

## HTML Source

```
1 <html>
2   <!--[...]-->
3   <h1>Individual papers</h1>
4   <h3>Invited Talks</h3>
5   <p>
6     <a id="[...]" class="citation_title"
          href="[...]">Data-Driven [...]</a>
7     <!--[...]-->
8     <span class="[...]title">9th [...]</span>
9     <span class="[...]firstpage">2</span>
10    <span class="[...]lastpage">2</span>
11    <span class="[...]pdf_url">http[...]</span>
12    <br/>
13    <span class="[...]author">Ra[...]</span>
14    <!--[...]-->
15  </p>
16  <!--[...]-->
17 </html>
```

## OXPath Expression

```
1 doc('http://edm2016.org/proceedings.html')
2   //p[./*[@class='cit']]:<record>
3     [./*[@class='author']:<authors=string(.)>]
4     [???:<title=???>]
5     [???:<pages=???>]
6     [???:<url=???>]
7     [???:<header=???>]
```

## XML Output

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <results>
3   <!--[...]-->
4   <record>
5     <authors>Rakesh Agrawal</authors>
6     <title>???</title>
7     <pages>???</pages>
8     <url>???</url>
9     <header>???</header>
10  </record>
11  <!--[...]-->
12 </results>
```

# EDM 2016: Adapting from EDM 2014

## HTML Source

```
 1 <html>
 2   <!--[...]-->
 3   <h1>Individual papers</h1>
 4   <h3>Invited Talks</h3>
 5   <p>
 6     <a id="[...]" class="citation_title"
          href="[...]">Data-Driven [...]</a>
 7     <!--[...]-->
 8     <span class="[...]title">9th [...]</span>
 9     <span class="[...]firstpage">2</span>
10     <span class="[...]lastpage">2</span>
11     <span class="[...]pdf_url">http[...]</span>
12     <br/>
13     <span class="[...]author">Ra[...]</span>
14     <!--[...]-->
15   </p>
16   <!--[...]-->
17 </html>
```

## OXPath Expression

```
1 doc('http://edm2016.org/proceedings.html')
2 //p[./*[@class='cit']]:<record>
3   [./*[@class='author']:<authors=string(.)>]
4   [./*[@class='title']:<title=string(.)>]
5   [???:<pages=???>]
6   [???:<url=???>]
7   [???:<header=???>]
```

## XML Output

```
 1 <?xml version="1.0" encoding="UTF-8"?>
 2 <results>
 3   <!--[...]-->
 4   <record>
 5     <authors>Rakesh Agrawal</authors>
 6     <title>Data-Driven [...]</title>
 7     <pages>???</pages>
 8     <url>???</url>
 9     <header>???</header>
10   </record>
11   <!--[...]-->
12 </results>
```

# EDM 2016: Adapting from EDM 2014

## HTML Source

```
 1 <html>
 2   <!--[...]-->
 3   <h1>Individual papers</h1>
 4   <h3>Invited Talks</h3>
 5   <p>
 6     <a id="[...]" class="citation_title"
        href="[...]">Data-Driven [...]</a>
 7     <!--[...]-->
 8     <span class="[...]title">9th [...]</span>
 9     <span class="[...]firstpage">2</span>
10     <span class="[...]lastpage">2</span>
11     <span class="[...]pdf_url">http[...]</span>
12     <br/>
13     <span class="[...]author">Ra[...]</span>
14     <!--[...]-->
15   </p>
16   <!--[...]-->
17 </html>
```

## OXPath Expression

```
1 doc('http://edm2016.org/proceedings.html')
2 //p[./*[@class-='cit']]:<record>
3   [./*[@class-='author']:<authors=string(.)>]
4   [./*[@class-='title']:<title=string(.)>]
5   [.:<pages=concat(./*[@class-='firstpage'],
        '-', ./*[@class-='lastpage'])>]
6   [???:<url=???>]
7   [???:<header=???>]
```

## XML Output

```
 1 <?xml version="1.0" encoding="UTF-8"?>
 2 <results>
 3   <!--[...]-->
 4   <record>
 5     <authors>Rakesh Agrawal</authors>
 6     <title>Data-Driven [...]</title>
 7     <pages>2-2</pages>
 8     <url>???</url>
 9     <header>???</header>
10   </record>
11   <!--[...]-->
12 </results>
```

# EDM 2016: Adapting from EDM 2014

## HTML Source

```
 1 <html>
 2   <!--[...]-->
 3   <h1>Individual papers</h1>
 4   <h3>Invited Talks</h3>
 5   <p>
 6     <a id="[...]" class="citation_title"
          href="[...]">Data-Driven [...]</a>
 7     <!--[...]-->
 8     <span class="[...]title">9th [...]</span>
 9     <span class="[...]firstpage">2</span>
10     <span class="[...]lastpage">2</span>
11     <span class="[...]pdf_url">http[...]</span>
12     <br/>
13     <span class="[...]author">Ra[...]</span>
14     <!--[...]-->
15   </p>
16   <!--[...]-->
17 </html>
```

## OXPath Expression

```
 1 doc('http://edm2016.org/proceedings.html')
 2   //p[./*[@class-='cit']]:<record>
 3     [./*[@class-='author']:<authors=string(.)>]
 4     [./*[@class-='title']:<title=string(.)>]
 5     [.:<pages=concat(./*[@class-='firstpage'],
             '-', ./*[@class-='lastpage'])>]
 6     [./*[@class-='url']:<url=string(.)>]
 7     [???:<header=???>]
```

## XML Output

```
 1 <?xml version="1.0" encoding="UTF-8"?>
 2 <results>
 3   <!--[...]-->
 4   <record>
 5     <authors>Rakesh Agrawal</authors>
 6     <title>Data-Driven [...]</title>
 7     <pages>2-2</pages>
 8     <url>http://[...].pdf</ee>
 9     <header>???</header>
10   </record>
11   <!--[...]-->
12 </results>
```

# EDM 2016: Adapting from EDM 2014

## HTML Source

```
 1 <html>
 2   <!--[...]-->
 3   <h1>Individual papers</h1>
 4   <h3>Invited Talks</h3>
 5   <p>
 6     <a id="[...]" class="citation_title"
          href="[...]">Data-Driven [...]</a>
 7     <!--[...]-->
 8     <span class="[...]title">9th [...]</span>
 9     <span class="[...]firstpage">2</span>
10     <span class="[...]lastpage">2</span>
11     <span class="[...]pdf_url">http[...]</span>
12     <br/>
13     <span class="[...]author">Ra[...]</span>
14     <!--[...]-->
15   </p>
16   <!--[...]-->
17 </html>
```

## OXPath Expression

```
1 doc('http://edm2016.org/proceedings.html')
2 //p[./*[@class='cit']]:<record>
3   [./*[@class-='author']:<authors=string(.)>]
4   [./*[@class-='title']:<title=string(.)>]
5   [.:<pages=concat(./*[@class-='firstpage'],
        '-', ./*[@class-='lastpage'])>]
6   [./*[@class-='url']:<url=string(.)>]
7   [./preceding::h3[1]:<header=string(.)>]
```

## XML Output

```
 1 <?xml version="1.0" encoding="UTF-8"?>
 2 <results>
 3   <!--[...]-->
 4   <record>
 5     <authors>Rakesh Agrawal</authors>
 6     <title>Data-Driven [...]</title>
 7     <pages>2-2</pages>
 8     <url>http://[...].pdf</ee>
 9     <header>Invited Talks</header>
10   </record>
11   <!--[...]-->
12 </results>
```

# Table of Contents

# Demonstration: Booking



## OXPath Expression

```
1 doc("https://www.booking.com/hotel/de/metropolitan...")
2  //a[@id='show_reviews_tab']/{click /}
3  /(//*[@id='review_next_page_link']/{clkwithchange
        /})*{0,2}
4  //div[contains(@class, 'review_list_block')]
5  //li[contains(@class, 'review_item')]
        [not(contains(@class, 'featured_review_item'))]
        [not(@class= 'review_item_photo
        review_item_photo-p')]:<review>
6  [? .//*[contains(@class, 'review_item_date')]
7  :<date=normalize-space(.)>]
8  [? .//*[contains(@class, 'review_item_review_score')]
9  :<score=normalize-space(.)>]
10 [? .//*[@class='review_item_header_content_container']
11 :<title=normalize-space(.)>]
12 [? .//*[contains(@class, 'review_item_review_content')]
13 :<text=string-join(./p/text(), " ")>]
14 [? .//*[@class='reviewer_country']
15 :<country=normalize-space(.)>]
```

# Demonstration: Twitter



## OXPath Expression

```
1 doc("https://twitter.com/search-home")
2 //input[@id='search-home-input']/{'urlaub'}/{pressenter/}
3 /(///div[contains(@class,'stream-footer')]
4   /{mouseover /})*{0, 4}
5 /.:<count=count(//li[@data-item-type='tweet'])>
6 //li[@data-item-type='tweet']:<tweet>
7   [? .//strong[@class='fullname show-popup-with-id ']
      :<user_name=string(.)>]
8   [? .//a[starts-with(@class,
      'account-group')]/span[@class='username u-dir']
      :<user_id=string(.)>]
9   [? .//a[@class="tweet-timestamp js-permalink js-nav
      js-tooltip"]/@title:<date=string(.)>]
10  [? .//p[starts-with(@class,'TweetTextSize')]
      :<content=normalize-space(.)>]
11  [? .//button[contains(@aria-describedby,
      'reply-count')]/span/span:<antworten=string(.)>]
12  [? .//button[contains(@aria-describedby,
      'retweet-count')]/span/span:<retweets=string(.)>]
13  [? .//button[contains(@aria-describedby,
      'favorite-count')]/span/span:<likes=string(.)>]
```