

**OPAL**  
OPEN DATA PORTAL

# Deliverable D3.3

## Erste Metadatenextraktionskomponente

Autoren: Adrian Wilke

Reviewer: Michael Röder

Veröffentlichung	Öffentlich
Fälligkeitsdatum	30. April 2019
Fertigstellung	7. Mai 2019
Arbeitspakete	A3.2 Semistrukturierte Metadatenextraktion A3.3 Unstrukturierte Metadatenextraktion
Typ	Software
Status	Final
Version	1.0

### Kurzfassung:

Die erste Version der OPAL Komponente zur Metadatenextraktion umfasst die Spracherkennung von Texten, die Eigennamenerkennung von Entitäten (Named Entity Recognition, NER) und die anschließende Disambiguierung. Außerdem werden Lizenzen für einzelne Distributionen von Metadaten extrahiert.

### Schlagworte:

Metadaten, Extraktion, Entitäten, Eigennamenerkennung, Disambiguierung, Spracherkennung, Lizenzen, FOX, AGDISTIS

# Inhalt



<a href="#"><u>Metadatenextraktion von Entitäten</u></a>	<a href="#"><u>2</u></a>
<a href="#"><u>Spracherkennung</u></a>	<a href="#"><u>2</u></a>
<a href="#"><u>Eigennamenerkennung und Disambiguierung von Entitäten</u></a>	<a href="#"><u>2</u></a>
<a href="#"><u>OPAL Webservices zur Metadatenextraktion</u></a>	<a href="#"><u>3</u></a>
<a href="#"><u>Anwendungsbeispiel</u></a>	<a href="#"><u>4</u></a>
<a href="#"><u>Spracherkennung</u></a>	<a href="#"><u>4</u></a>
<a href="#"><u>Eigennamenerkennung und Disambiguierung von Entitäten</u></a>	<a href="#"><u>4</u></a>
<a href="#"><u>Metadatenextraktion von Lizenzen</u></a>	<a href="#"><u>5</u></a>
<a href="#"><u>Datenschema: Einordnung von Lizenzinformationen</u></a>	<a href="#"><u>5</u></a>
<a href="#"><u>Übersicht extrahierter Lizenz-URIs</u></a>	<a href="#"><u>6</u></a>

### 1 Metadatenextraktion von Entitäten

Die Extraktion von Entitäten ist dem Arbeitspaket *A3.3 Unstrukturierte Metadatenextraktion* zugeordnet. Die zugehörige OPAL Softwarekomponente zu diesem Arbeitspaket wurde als OPAL Metadata-Extraction<sup>1</sup> veröffentlicht.

Zur Anreicherung von Daten wird eine unterliegende Datenbasis analysiert und die resultierenden Ergebnisse dem ursprünglichen Datenmodell zurückgespiegelt. Die unterliegende Datenbasis in OPAL setzt sich mehrheitlich aus Metadaten offener Datensätze zusammen. Dabei wird für einzelne Bestandteile der Metadaten natürlichsprachlicher Text verwendet. Dies ist beispielsweise in Beschreibungstexten der Fall, die den Inhalt der referenzierten Datensätze zusammenfassen. Diese Beschreibungstexte enthalten häufig Entitäten, z.B. Städtenamen. Befinden sich solche Entitäten in weiteren Datensätzen, so können sie verwendet werden, um Benutzern eines Informationssystems Vorschläge zu weiteren Datensätzen zu unterbreiten, die ebenfalls interessant sein könnten. Hat ein Benutzer beispielsweise einen Datensatz zu einer bestimmten Stadt gefunden, interessiert er sich möglicherweise für weitere Datensätze zur selben Stadt. Die entwickelte Komponente zur Erkennung und Extraktion von Entitäten kann also als Grundlage dazu dienen, Vorschläge für Datensätze zu generieren, die die Auffindbarkeit von Daten verbessern.

Die Beschreibung der Komponente zur Metadatenextraktion setzt sich aus vier Teilen zusammen. Die (1) Spracherkennung dient dazu, ein Erkennungs-Modell auszuwählen, das für die jeweiligen Sprache erstellt wurde. Anschließend werden (2) Entitäten gesucht und eindeutige Bezeichner für diese zurückgegeben. Diese beiden Komponenten können über (3) Webservices zugegriffen werden. Hierzu wird ein (4) Anwendungsbeispiel angeführt.

#### 1.1 Spracherkennung

Zur Erkennung, in welcher Sprache ein Text verfasst wurde, ist in der OPAL Metadatenextraktion die Software Apache OpenNLP<sup>2</sup> eingebunden. Im Zuge dessen wird das Language Detector Model<sup>3</sup> in Version 1.8.3 eingesetzt. Damit können 103 Sprachen erkannt werden. Diese werden im Modell durch das ISO 639-3<sup>4</sup> Format zugeordnet. In OPAL wird übergreifend das ISO 639-1 Format verwendet, welches im Data Catalog Vocabulary DCAT<sup>5</sup> empfohlen wird. Daher erfolgt intern eine Übersetzung zwischen den beiden ISO 639 Standards.

#### 1.2 Eigennamenerkennung und Disambiguierung von Entitäten

Die in OPAL fokussierten Metadaten beschreiben Mobilitätsdaten. Diese beziehen sich häufig auf eine Region, z.B. eine Stadt. Die Metadaten enthalten als Beschreibung unstrukturierte, natürlichsprachliche Texte, die wiederum Entitäten beinhalten. Beispiele für Entitäten sind Städte oder Bundesländer, für die weitere Daten (z.B. Geo-Koordinaten) abgerufen werden können. Zur Erkennung von Eigennamen (z.B. Paderborn) wird in OPAL die Software FOX<sup>6</sup> eingesetzt. Die anschließende Disambiguierung (die Erkennung von eindeutigen Bezeichnern) von Entitäten (z.B. <http://de.dbpedia.org/resource/Paderborn>) geschieht über die Einbindung der Software AGDISTIS<sup>7</sup>. Die OPAL Komponenten zur Metadatenextraktion sind über Webservices aufrufbar.

---

<sup>1</sup> OPAL Metadata-Extraction: <https://github.com/projekt-opal/metadata-extraction>

<sup>2</sup> Apache OpenNLP: <https://opennlp.apache.org/>

<sup>3</sup> Language Detector Model: <https://opennlp.apache.org/models.html>

<sup>4</sup> ISO 639-3: <https://iso639-3.sil.org/>

<sup>5</sup> ISO 639-1 Verweis in DCAT: [https://www.w3.org/TR/vocab-dcat/#Property:catalog\\_language](https://www.w3.org/TR/vocab-dcat/#Property:catalog_language)

<sup>6</sup> FOX - Federated knOwledge eXtraction Framework: <https://dice.cs.upb.de/projects/active-projects/fox/>

<sup>7</sup> AGDISTIS - Agnostic Named Entity Disambiguation: <https://github.com/dice-group/AGDISTIS>

### 1.3 OPAL Webservices zur Metadatenextraktion

In der folgenden Tabelle sind die entwickelten Webservices zur Metadatenextraktion beschrieben. Dies umfasst eine Kurzbeschreibung der Webservices, die benötigten Eingabeparameter und jeweils ein Beispiel zum Aufruf der Webservices.

Spracherkennung	
Beschreibung:	Gibt die erkannte Sprache im Format ISO 639-3 zurück.
Parameter:	Volltext
Beispiel:	<a href="http://opal/metadata/lang?text=Sprachen lernen">http://opal/metadata/lang?text=Sprachen lernen</a>
Beschreibung:	Erkennt die verwendeten Sprachen und fügt sie den OPAL RDF-Datensätzen hinzu.
Parameter:	RDF-Modell im TURTLE-Format
Beispiel:	<a href="http://opal/metadata/lang/model?turtleBytes=...">http://opal/metadata/lang/model?turtleBytes=...</a>
Eigennamenerkennung und Disambiguierung von Entitäten	
Beschreibung:	Gibt das FOX Ergebnis im Format RDF-TURTLE zurück.
Parameter:	Volltext, Sprache
Beispiel:	<a href="http://opal/metadata/fox?text=A. Einstein was born in Ulm.&amp;lang=en">http://opal/metadata/fox?text=A. Einstein was born in Ulm.&amp;lang=en</a>
Beschreibung:	Gibt die natürlichsprachlichen Namen von Orten im JSON Format zurück.
Parameter:	Volltext, Sprache
Beispiel:	<a href="http://opal/metadata/fox/location/names?text=Paderborn und Bad Oeynhausen sind in NRW&amp;lang=de">http://opal/metadata/fox/location/names?text=Paderborn und Bad Oeynhausen sind in NRW&amp;lang=de</a>
Beschreibung:	Gibt die URIs von Orten im JSON Format zurück.
Parameter:	Volltext, Sprache
Beispiel:	<a href="http://opal/metadata/fox/location/uris?text=Paderborn und Bad Oeynhausen sind in NRW&amp;lang=de">http://opal/metadata/fox/location/uris?text=Paderborn und Bad Oeynhausen sind in NRW&amp;lang=de</a>

**Tabelle: Entwickelte OPAL Webservices zur Metadatenextraktion**

### 1.4 Anwendungsbeispiel

Ein Anwendungsbeispiel zur Metadatenextraktion ist der folgende Beispieltext, der stellvertretend für eine Beschreibung eines Datensatzes verwendet wird.

Beispieltext:

```
Beispielanfrage zu Paderborn und Bad Oeynhausen in NRW
```

#### 1.4.1 Spracherkennung

Für den Beispieltext wird zunächst die verwendete Sprache bestimmt.

Aufruf Webservice Spracherkennung:

```
http://opal/metadata/lang?text=Beispielanfrage%20zu%20Paderborn%20und%20Bad%20Oeynhaus  
sen%20in%20NRW
```

Resultat:

```
deu
```

Die Rückgabe ist ein 3-stelliger CODE im ISO 639-3 Format und repräsentiert die Sprache Deutsch.

#### 1.4.2 Eigennamenerkennung und Disambiguierung von Entitäten

Die Metadatenextraktion der Geo-Entitäten geschieht über den folgenden Aufruf.

Webservice:

```
http://opal/metadata/fox/location/names?text=Beispielanfrage%20Paderborn%20und%20Bad%  
20Oeynhausen%20in%20NRW&lang=de
```

Resultat:

```
[ "Bad Oeynhausen", "NRW", "Paderborn" ]
```

In diesem Beispiel werden eine Stadt bestehend aus einem Wort, eine Stadt aus zwei Wörtern und die Abkürzung eines Bundeslandes erkannt. Eine weitere Variante mit Disambiguierung zeigt der folgende Aufruf:

Webservice:

```
http://opal/metadata/fox/location/uris?text=Beispielanfrage%20Paderborn%20und%20Bad%  
20Oeynhausen%20in%20NRW&lang=de
```

Resultat:

```
[ "http://de.dbpedia.org/resource/Bad_Oeynhausen",  
"http://de.dbpedia.org/resource/NRW", "http://de.dbpedia.org/resource/Paderborn" ]
```

Bei diesem Aufruf werden URIs der Entitäten zurückgegeben. Diese sind eindeutige Bezeichner für die jeweiligen Entitäten. Dadurch können sie weiterverwendet werden, um z.B. entsprechende Geokoordinaten zu ermitteln. Auf der DBpedia Seite zur Ressource Paderborn<sup>8</sup> findet sich eine Übersicht abrufbarer Metadaten.

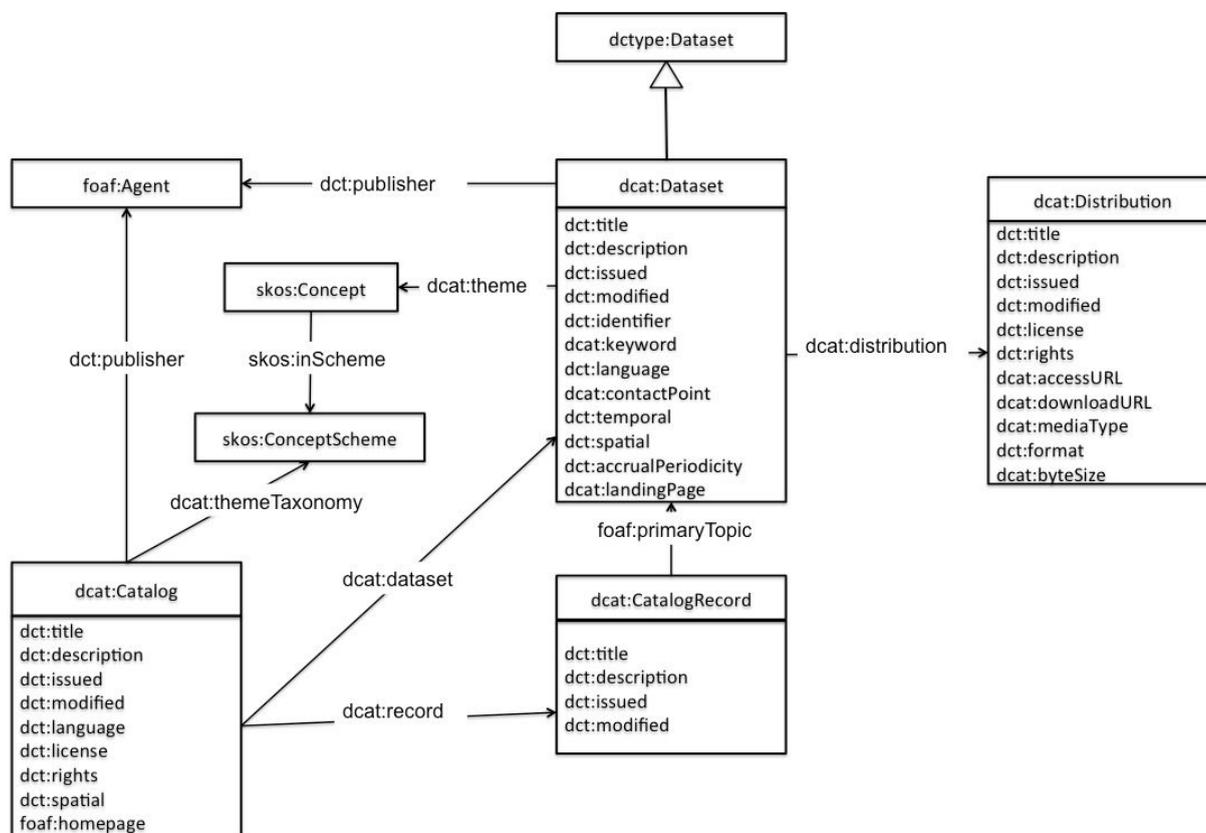
<sup>8</sup> DBpedia Ressource Paderborn: <http://dbpedia.org/page/Paderborn>

## 2 Metadatenextraktion von Lizenzen

Die Extraktion von Lizenzinformationen ist dem Arbeitspaket *A3.2 Semistrukturierte Metadatenextraktion* zugeordnet.

### 2.1 Datenschema: Einordnung von Lizenzinformationen

Der Extraktion von Metadaten aus offenen Datenportalen, die für OPAL genutzt werden, schließt sich eine Datenspeicherung im RDF-Format an. Hierzu wird u.a. das DCAT Vokabular<sup>9</sup> genutzt. Datenportale entsprechen dabei dem Typ `dcat:Catalog`. Die dort enthaltenen Metadatenätze werden in OPAL als `dcat:Dataset` abgebildet. Diese Datensätze können u.a. in verschiedenen Datenformaten (z.B. PDF, XLSX, CSV, RDF-TURTLE) zugänglich gemacht werden. Solche Daten-Versionen werden als Distribution bezeichnet und als `dcat:Distribution` klassifiziert. Für Distributionen kann über den Bezeichner `dct:license` eine URI angegeben werden, die auf die jeweils verwendete Lizenz verweist. Die folgende Abbildung veranschaulicht das verwendete DCAT Vokabular:



**Abbildung: DCAT Vokabular**  
(Quelle: <https://www.w3.org/TR/vocab-dcat/>)

<sup>9</sup> DCAT: <https://www.w3.org/TR/vocab-dcat/>

## 2.2 Übersicht extrahierter Lizenz-URIs

Für die Datenportale mCLOUD<sup>10</sup> und European Data Portal<sup>11</sup> wurden Lizenz-URIs für 2.798 und 61.173 einzelne Distributionen extrahiert. Die Analyse der Lizenzen geschah über Abfragen in der Sprache SPARQL und entsprechendem Programmcode<sup>12</sup>. Die folgende Tabelle zeigt eine Übersicht der verwendeten Lizenzen:

Lizenz URI	Anzahl Distributionen
<a href="http://europeandataportal.eu/content/show-license?license_id=OGL2.0">http://europeandataportal.eu/content/show-license?license_id=OGL2.0</a>	34.911
<a href="https://www.europeandataportal.eu/content/show-license?license_id=OGL2.0">https://www.europeandataportal.eu/content/show-license?license_id=OGL2.0</a>	10.645
<a href="http://europeandataportal.eu/content/show-license?license_id=FR-LO">http://europeandataportal.eu/content/show-license?license_id=FR-LO</a>	9.493
<a href="https://www.europeandataportal.eu/content/show-license?license_id=ODC-ODbL">https://www.europeandataportal.eu/content/show-license?license_id=ODC-ODbL</a>	2.866
<a href="http://europeandataportal.eu/content/show-license?license_id=ODC-ODbL">http://europeandataportal.eu/content/show-license?license_id=ODC-ODbL</a>	1.420
<a href="https://www.europeandataportal.eu/content/show-license?license_id=FR-LO">https://www.europeandataportal.eu/content/show-license?license_id=FR-LO</a>	1.293
<a href="file:///content/show-license?license_id=OGL2.0">file:///content/show-license?license_id=OGL2.0</a>	870
<a href="http://www.gesetze-im-internet.de/geonutzv/index.html">http://www.gesetze-im-internet.de/geonutzv/index.html</a>	758
<a href="https://www.govdata.de/dl-de/by-2-0">https://www.govdata.de/dl-de/by-2-0</a>	671
<a href="https://creativecommons.org/licenses/by/3.0/de/">https://creativecommons.org/licenses/by/3.0/de/</a>	562
<a href="https://creativecommons.org/publicdomain/zero/1.0/deed.de">https://creativecommons.org/publicdomain/zero/1.0/deed.de</a>	244
<a href="https://creativecommons.org/licenses/by/4.0/deed.de">https://creativecommons.org/licenses/by/4.0/deed.de</a>	193
<a href="http://www.opendefinition.org/licenses/cc-by">http://www.opendefinition.org/licenses/cc-by</a>	140
<a href="http://www.stadtentwicklung.berlin.de/geoinformation/download/nutzIII.pdf">http://www.stadtentwicklung.berlin.de/geoinformation/download/nutzIII.pdf</a>	94
<a href="http://europeandataportal.eu/content/show-license?license_id=DL-DE-BY2.0">http://europeandataportal.eu/content/show-license?license_id=DL-DE-BY2.0</a>	76
<a href="https://www.govdata.de/dl-de/zero-2-0">https://www.govdata.de/dl-de/zero-2-0</a>	75
<a href="file:///content/show-license?license_id=FR-LO">file:///content/show-license?license_id=FR-LO</a>	67
<a href="file:///content/show-license?license_id=ODC-ODbL">file:///content/show-license?license_id=ODC-ODbL</a>	41
<a href="file:///content/show-license?license_id=CC-BY4.0">file:///content/show-license?license_id=CC-BY4.0</a>	40
<a href="http://www.ioer-monitor.de/fileadmin/Dokumente/PDFs/Nutzungsbedingungen_IOER-Monitor.pdf">http://www.ioer-monitor.de/fileadmin/Dokumente/PDFs/Nutzungsbedingungen_IOER-Monitor.pdf</a>	30
<a href="file:///content/show-license?license_id=DL-DE-BY2.0">file:///content/show-license?license_id=DL-DE-BY2.0</a>	16
<a href="https://scihub.copernicus.eu/twiki/pub/SciHubWebPortal/TermsConditions/TC_Sentinel_Data_31072014.pdf">https://scihub.copernicus.eu/twiki/pub/SciHubWebPortal/TermsConditions/TC_Sentinel_Data_31072014.pdf</a>	16
<a href="http://europeandataportal.eu/content/show-license?license_id=GFDL-1.3">http://europeandataportal.eu/content/show-license?license_id=GFDL-1.3</a>	13
<a href="https://www.europeandataportal.eu/content/show-license?license_id=CC-BY4.0">https://www.europeandataportal.eu/content/show-license?license_id=CC-BY4.0</a>	12
<a href="http://europeandataportal.eu/content/show-license?license_id=ODC-BY">http://europeandataportal.eu/content/show-license?license_id=ODC-BY</a>	7
<a href="http://europeandataportal.eu/content/show-license?license_id=CC-BY4.0">http://europeandataportal.eu/content/show-license?license_id=CC-BY4.0</a>	6

<sup>10</sup> mCLOUD: <https://www.mcloud.de/>

<sup>11</sup> European Data Portal: <https://www.europeandataportal.eu/>

<sup>12</sup> Metadatenextraktion für Statistiken: <https://github.com/projekt-opal/statistics>

### D3.3 - Erste Metadatenextraktionskomponente

<a href="http://opendatacommons.org/licenses/odbl/summary/">http://opendatacommons.org/licenses/odbl/summary/</a>	6
<a href="file:///content/show-license?license_id=DL-DE-BY1.0">file:///content/show-license?license_id=DL-DE-BY1.0</a>	5
<a href="file:///content/show-license?license_id=ODC-BY">file:///content/show-license?license_id=ODC-BY</a>	4
<a href="http://www.opendefinition.org/licenses/cc-zero">http://www.opendefinition.org/licenses/cc-zero</a>	3
<a href="http://www.opendefinition.org/licenses/odc-odbl">http://www.opendefinition.org/licenses/odc-odbl</a>	3
<a href="https://www.europeandataportal.eu/content/show-license?license_id=ODC-BY">https://www.europeandataportal.eu/content/show-license?license_id=ODC-BY</a>	3
<a href="file:///content/show-license?license_id=GFDL-1.3">file:///content/show-license?license_id=GFDL-1.3</a>	2
<a href="http://europeandataportal.eu/content/show-license?license_id=PSEUL">http://europeandataportal.eu/content/show-license?license_id=PSEUL</a>	2
<a href="http://images.vbb.de/assets/downloads/file/20928.pdf">http://images.vbb.de/assets/downloads/file/20928.pdf</a>	1
<a href="http://www.opendefinition.org/licenses/cc-by-sa">http://www.opendefinition.org/licenses/cc-by-sa</a>	1
<a href="https://www.europeandataportal.eu/content/show-license?license_id=GFDL-1.3">https://www.europeandataportal.eu/content/show-license?license_id=GFDL-1.3</a>	1
<a href="https://www.vrsinfo.de/fileadmin/Dateien/api/NutzervereinbarungODOS.pdf">https://www.vrsinfo.de/fileadmin/Dateien/api/NutzervereinbarungODOS.pdf</a>	1
Gesamt	64.591

**Tabelle: Übersicht extrahierter Lizenzen**

Für diese extrahierten Lizenz-URIs besteht in den Arbeitspaketen im weiteren Projektverlauf die Möglichkeit, eine erweiterte Datenaufbereitung durchzuführen. Es bietet sich eine Disambiguierung der Lizenzen und der verwendeten URIs an. Die DCAT-Erweiterung DCAT-AP.de<sup>13</sup> bietet in Version 1.0 eine Klassifizierung von 33 Lizenzen an, in denen angegeben ist, ob der jeweils verwendete Lizenztyp offen oder geschlossen ist.

DCAT-AP.de wurde beim 4. Treffen des WIK/mFUND Arbeitsforums Standardisierung/mCLOUD<sup>14</sup> am 20. März 2019 in der Präsentation zum Datenportal GovData vorgestellt. Die Veranstaltung fand beim Deutschen Zentrum für Luft- und Raumfahrt (DLR) in Bremen und mit Unterstützung vom Bundesministerium für Verkehr und digitale Infrastruktur (BMVI) statt.

Weitere Dokumente zum Projekt OPAL sind auf der Webseite verfügbar:

<http://projekt-opal.de/en/results/deliverables/>

<sup>13</sup> DCAT-AP.de: <https://www.dcat-ap.de/def/>

<sup>14</sup> WIK/mFUND Arbeitsforum Standardisierung/mCLOUD: <https://www.wik.org/index.php?id=943>