

# Semantische Suche

Autor(en): **Dengel, Andreas**

Objektyp: **Article**

Zeitschrift: **Bulletin.ch : Fachzeitschrift und Verbandsinformationen von Electrosuisse, VSE = revue spécialisée et informations des associations Electrosuisse, AES**

Band (Jahr): **103 (2012)**

Heft 4

PDF erstellt am: **12.07.2024**

Persistenter Link: <https://doi.org/10.5169/seals-857287>

## **Nutzungsbedingungen**

Die ETH-Bibliothek ist Anbieterin der digitalisierten Zeitschriften. Sie besitzt keine Urheberrechte an den Inhalten der Zeitschriften. Die Rechte liegen in der Regel bei den Herausgebern.

Die auf der Plattform e-periodica veröffentlichten Dokumente stehen für nicht-kommerzielle Zwecke in Lehre und Forschung sowie für die private Nutzung frei zur Verfügung. Einzelne Dateien oder Ausdrucke aus diesem Angebot können zusammen mit diesen Nutzungsbedingungen und den korrekten Herkunftsbezeichnungen weitergegeben werden.

Das Veröffentlichen von Bildern in Print- und Online-Publikationen ist nur mit vorheriger Genehmigung der Rechteinhaber erlaubt. Die systematische Speicherung von Teilen des elektronischen Angebots auf anderen Servern bedarf ebenfalls des schriftlichen Einverständnisses der Rechteinhaber.

## **Haftungsausschluss**

Alle Angaben erfolgen ohne Gewähr für Vollständigkeit oder Richtigkeit. Es wird keine Haftung übernommen für Schäden durch die Verwendung von Informationen aus diesem Online-Angebot oder durch das Fehlen von Informationen. Dies gilt auch für Inhalte Dritter, die über dieses Angebot zugänglich sind.

# Semantische Suche

## Auf dem Weg zur Bedeutungserschliessung grosser Informationsmengen

Internet-Suchmaschinen haben ihren Ursprung im Information Retrieval. Sie erstellen einen Index aus Schlüsselworten, der für Suchanfragen mittels syntaktischem Vergleich verwendet wird, und liefern eine nach Relevanz geordnete Trefferliste. Semantische Suchsysteme hingegen erkennen automatisch den Bedeutungszusammenhang eines Suchterms und schlagen alternative Begriffe vor. Der semantische Ansatz verspricht, die Suche besser, intuitiver und relevanter zu gestalten.

### Andreas Dengel

Gibt man in eine konventionelle Suchmaschine den Begriff «Bank» ein, so erhält man einige Hundert Millionen Treffer in Form von Verweisen auf möglicherweise relevante Dokumente, die den Suchbegriff enthalten. Da die traditionelle Schlüsselwortsuche lediglich das Vorkommen der Suchbegriffe im Index überprüft, bleiben Akronyme oder Synonyme mit gleicher oder ähnlicher Bedeutung unberücksichtigt. Für die Anfrage nach «Bank» werden relevante Dokumente, die beispielsweise den Begriff «Geldinstitut» enthalten, nicht gefunden. Semantische Suchsysteme erlauben in solchen Fällen beispielsweise Rückfragen wie «Meinten Sie das Finanzinstitut oder

die Sitzgelegenheit?». Dazu werden Bedeutungen der Anfrage, wie auch Aussagen der zu durchsuchenden Texte mithilfe kognitiver Algorithmen analysiert und es wird versucht, diese formal zu verstehen. Man kann die semantische Suche als einen Suchprozess betrachten, bei dem in jeder Phase der Suche semantisches Wissen verwendet wird. Dies betrifft die Anfragestellung genauso wie die (eigentliche) Suche und Ergebnisrepräsentation.

### Grenzen klassischer Suche

Nehmen wir an, wir möchten die Adresse des Financial Service Unternehmens CarFS in Bern herausfinden und geben

die Begriffe «Adresse», «CarFS», «Bern» in die Suchmaschine Google ein, so erhalten wir alle Dokumente, in denen die genannten Stichworte vorkommen. Da die Begriffe als Zeichenketten betrachtet werden, versteht die Suchmaschine weder, versteht die Suchmaschine weder, worum es geht, noch, dass die Anfrage ein ganz bestimmtes Faktum anfragt. Daher umfasst die Ergebnisliste auch ca. 71 Millionen potenziell relevanter Dokumente, mit Vorschlägen wie «Einkaufen in Bern», Autovermietungen, Busreisen, TicketserVICES, usw. Wüsste die Suchmaschine, dass Bern eine Stadt in der Schweiz ist, so könnte sie einen Zusammenhang herstellen, dass es wohl um eine Adresse in Bern geht und hinter CarFS eine Person oder ein Unternehmen steht, das einen solchen Namen trägt. Entsprechend eingeschränkt kann dann eine gezielte Faktensuche stattfinden.

### Formularbasierte Suche

Eine Form zur näheren Bestimmung der Bedeutung wird mit dem Begriff «Formularbasierte Suche» umschrieben. Sie beruht darauf, dass ein Benutzer seine Anfrage unter Verwendung von (Web-) Formularen eingeben kann. Das Formular umfasst eine Menge von bezeichneten Feldern, wie «Stadt», in die man einen dazugehörigen Wert wie «Bern» eingeben kann. Da mit der Eingabe die Semantik

The screenshot shows the SIG.MA Semantic Information Mashup interface. The search query is 'dengel Andreas'. The results are displayed in a list on the right side of the page. The first result is 'Andreas Dengel' with 11 facts, dated 2009-09-03. The second result is 'Andreas Dengel' with 15 facts, dated 2010-11-24. The third result is 'Andreas Dengel - semanti...' with 3 facts, dated 2010-11-24. The fourth result is 'Andreas Dengel: books by...' with 69 facts, dated 2010-04-15. The fifth result is 'Andreas Dengel - Linked...' with 10 facts, dated 2009-10-30. The sixth result is 'Andreas Dengel - Compute...' with 3 facts, dated 2009-10-30. The seventh result is 'Andreas Dengel' with 11 facts, dated 2010-01-14. The eighth result is 'Andreas Dengel' with 9 facts, dated 2008-12-15. The ninth result is 'Andreas Dengel' with 9 facts, dated 2008-12-16. The tenth result is 'Andreas Dengel' with 9 facts, dated 2008-12-16. The eleventh result is 'Andreas Dengel' with 9 facts, dated 2008-12-16. The twelfth result is 'Andreas Dengel' with 9 facts, dated 2008-12-16. The thirteenth result is 'Andreas Dengel' with 9 facts, dated 2008-12-16.

Bild 1 Semantikbasierte Schlüsselwortsuche mit SIG.MA.



bereits festgelegt ist, sind für diese Form der Suche keine aufwendigen Suchalgorithmen erforderlich, denn die gemachten Angaben reichen aus, um die Abfrage, z.B. als SPARQL-Ausdruck, zu konstruieren und auszuführen.

### Klassifikationssysteme beschreiben Vokabular

An diesem Beispiel kann man bereits eine wichtige Dimension von Wissen erkennen, auf der semantische Suchsysteme aufbauen: Sie nutzen das Konzept der Klassifikation, das es erlaubt, gleichartige Ressourcen unter einem Begriff zusammenzufassen. Klassen sind dabei wie Variablen zu betrachten, d.h. sie können eine Menge von Werten annehmen, so wie in unserem Beispiel es neben Bern auch noch viele andere Städte gibt, wie Zürich, Basel oder Genf. Sie alle gehören damit zur Klasse «Stadt». Wenn man es genauer nimmt, so gehören sie sogar zur Klasse «Schweizer Stadt», die nur Werte von Städten annehmen kann, die auch in der Schweiz liegen. Man kann also sagen, dass «Schweizer Stadt» eine Unterklasse von «Stadt» ist.

An diesem erweiterten Beispiel wird nun noch klarer, welche Prinzipien einer formalen Semantik unterliegen. Eine Domäne (Vokabular) wird in Klassen (Begrifflichkeiten) gegliedert, diese werden in Form einer Klassenhierarchie (Taxonomie) angeordnet und am Ende werden jeder Klasse Werte zugeordnet. Diese Werte bezeichnet man als Instanzen einer Klasse (Bern ist eine Instanz von Stadt). Beides zusammen bildet eine Ontologie, die, formal beschrieben, das Hintergrundwissen für semantische Suchmaschinen darstellt.

### Ontologien als formale Semantik

Im Internet stehen bereits eine ganze Reihe solcher Ontologien als Linked Open Data [1] zur Verfügung. Ein prominenter Vertreter ist DBpedia – die semantische Version der offenen, freien Online-Enzyklopädie Wikipedia – wo die als Datensammlung vorliegenden HTML-Beschreibungen aus Wikipedia in RDF transformiert wurden. Das Resource Description Framework (RDF) [1] ist eine Standardsprache des W3C (World Wide Web Consortium) und dient der Beschreibung von Ontologien. Das zugrundeliegende Modell besteht aus den drei Objekttypen: Ressourcen, Eigenschaftselementen und Objekten. Jeweils eine Ressource, eine Eigenschaft

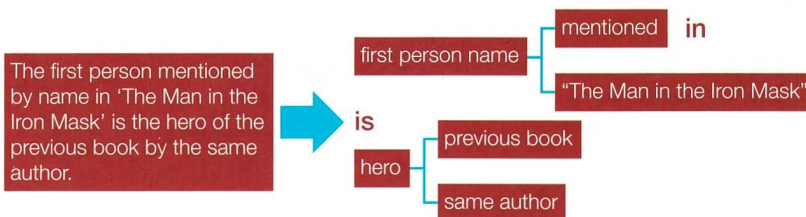


Bild 2 Anfragezerlegung bei Watson [6].

und ein Objekt bilden zusammen ein sogenanntes RDF-Tripel. Ein solches Tripel kann man sich wie einen elementaren Satz bestehend aus Subjekt, Prädikat und Objekt vorstellen, der eine Aussage über ein bestimmtes Objekt innerhalb einer Domäne macht. Die Aussage (das Tripel) etwa, dass Bern eine Stadt ist, lässt sich unter Zuhilfenahme des in DBpedia standardisierten und auf URIs (Uniform Resource Identifier) aufbauenden Vokabulars in RDF wie folgt beschreiben:

```
<http://dbpedia.org/resource/Bern>
<http://www.w3.org/1999/02/22-rdf-syntax-ns#type>
<http://dbpedia.org/ontology/City>
```

Eine aktuelle semantische Beschreibung für Bern findet man beispielsweise unter der Web-Adresse <http://dbpedia.org/page/Bern>.

### Formale Anfrageformulierung

Man kann die Nützlichkeit solcher Ontologien sehr schön am Beispiel von Suchmaschinen erkennen, deren Schwerpunkt in der Anfrageformulierung liegt. Entsprechende Suchmaschinen verwenden eine spezielle Sprache zur Formulierung von Anfragen, die sich der formalen Beschreibung von RDF bedient. SPARQL [2] beispielsweise ist eine solche formale Anfragesprache, mit der sich Suchanfragen der folgenden Form ausdrücken lassen:

```
PREFIX abc: <http://example.com/example-Ontologie#>
SELECT ?canton_capital ?canton
WHERE {
  ?x abc:cityname ?canton_capital.
  ?y abc:statename ?canton.
  ?x abc:isCapitalOf ?y.
  ?y abc:isInCountry abc:switzerland.
}
```

Im dargestellten Beispiel werden alle Kantonshauptstädte mit dem zugehörigen Schweizer Kanton ermittelt, sofern diese

in der Hintergrundontologie enthalten sind. Hierzu erfolgt eine Abfrage mittels Zugriff auf die aus RDF bekannten Tripel. Dabei werden alle Variablenbelegungen für `?canton_capital` und `?canton` zurückgegeben, die auf die Muster dieser vier RDF-Tripel passen. Durch den verwendeten Präfix kann die Notation verkürzt und somit leserlicher gestaltet werden.

### Semantikbasierte Schlüsselwortsuche

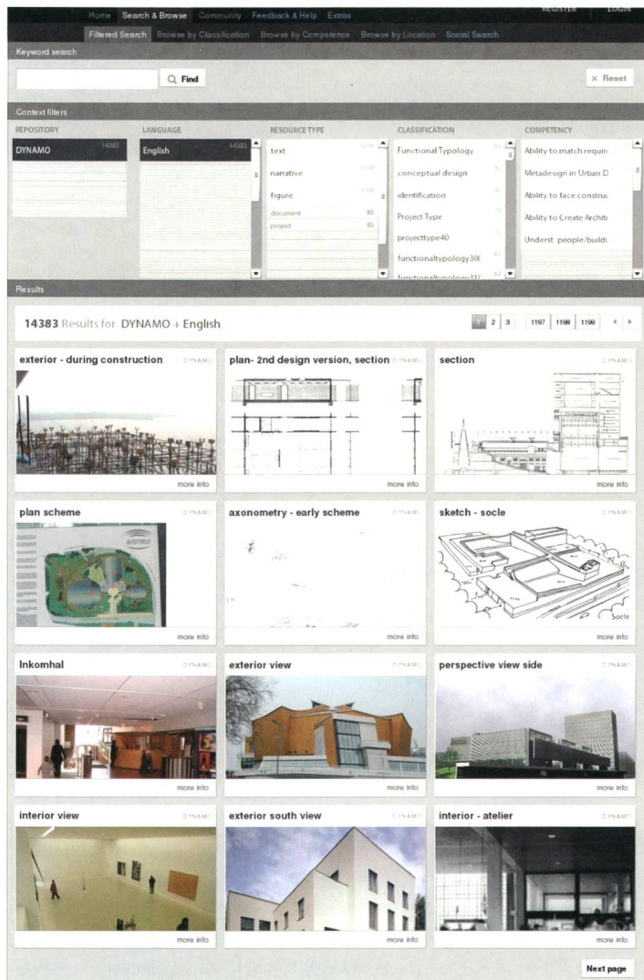
Aber auch reine Schlüsselwortsuche kann durch Semantik erweitert werden. Die ersten Suchmaschinen, die dieses Prinzip eingesetzt haben, verwendeten eine einfache Syntax, mit der man beschreiben konnte, welcher Klasse die Wörter einer Suchanfrage zugehörig sind. Durch erweiterte Anfragen wie `organization:bank city:zuerich` konnte man die Bedeutung der Suchterme präzisieren.

Ein aktuelles Beispiel für semantikbasierte Anfragesysteme ist die Suchmaschine SIG.MA (Semantic Information Mashup – <http://sig.ma/>). Sie durchstöbert in Webseiten eingebettete RDF-Annotationen, aggregiert die gefundenen Ressourcen und stellt diese strukturiert zur Verfügung. Gibt man etwa den Namen des Autors dieses Artikels in SIG.MA ein, so erhält man – wie in Bild 1 ausschnittsweise dargestellt – eine Menge von Fakten (links) und dazu gehörigen Ressourcen (rechts), die über die entsprechende URI zugreifbar sind. Bei der Eingabe kann der Benutzer direkt mitverfolgen, wie sukzessive immer mehr Fakten gefunden und geprüft werden und die Ergebnisliste Schritt für Schritt erweitert wird. Bemerkenswert für SIG.MA ist, dass es nicht nur Fakten hinzufügt, sondern während der Suche auch solche Fakten wieder entfernt, deren Mehrdeutigkeiten aufgelöst werden konnten.

### Question Answering

Eine erweiterte Form der semantischen Suche sind sogenannte Question-





**Bild 3** «Click & Filter» bei facetierter Suche.

Answering-Systeme. Sie sind dafür geeignet, sprachlich ausformulierte natürliche Anfragen mithilfe einer vorher definierten Semantik zu verarbeiten und zu beantworten. Sie verwenden oft auch weitergehendes linguistisches Wissen, um eine Anfrage zu analysieren und in eine formale Anfrage zu übersetzen. Ebenso werden die gefundenen, formal beschriebenen Ergebnisse in natürliche Sprache umgeformt, bevor sie als Antwortsatz ausgegeben werden. Die Wissensbasis der Question-Answering-Tools ist aufwendig konstruiert und bietet detailliertes Wissen über die vorhandenen Instanzen und deren Beziehungen. Ein bekannter Vertreter von Question-Answering ist das von IBM entwickelte System Watson [3], das englische Fragen – wie in **Bild 4** dargestellt – versteht, die enthaltenen Wörter im Kontext der Anfrage einordnet und mit Hilfe einer umfangreichen Wissensbasis schnell und präzise beantwortet.

Da semantische Suchsysteme auf umfangreichem Wissen aufbauen, wäre es natürlich sinnvoll, wenn der Benutzer für die Fragestellung das Vokabular der verwendeten Ontologien kennen würde. In

vielen Fällen werden zur Unterstützung entsprechende grafische Benutzerschnittstellen angeboten, die z.B. vorhandene Klassen, wie Orte (Länder, Städte, Gemeinden), Ereignisse (Feste, Kongresse, Aufführungen), Personen oder Organisationen z.B. mittels Pull-down-Menü zur Verfügung stellen, sowie zugehörige Prädikate (Eigenschaftselemente), z.B. LivesIn, WorksFor, oder HasLocation anzeigen, mit der Aussagen über die dazugehörigen Instanzen gemacht werden können.

### Facettiertes Browsing

Eine einfache, aber wirkungsvolle Variante dieser Form der semantischen Suche ist das sogenannte facetierte Browsing, das eine klassische Suche mit Klassifikationssystemen verbindet. Neben der bekannten Eingabe von Stichwörtern kann der Benutzer aus einer Reihe von gruppierten Klassen auswählen, die ihm helfen sollen, das zugrunde liegende Vokabular besser zu verstehen und seinen Informationsbedarf auf Grundlage des Informationsangebotes gezielter zu beschreiben bzw. zu nutzen. Solche polyhierarchischen Klassifikationssysteme

nennt man auch Facetten. Facetten beschreiben i.d.R. das Vokabular unterschiedlicher Informationssichten, wie z.B. Was, Wer, Wo, usw. und erlauben durch Anklicken der entsprechenden Klasse in den Facetten dann Anfragen wie etwa «Privatdarlehen, Banken, Zürich» sehr intuitiv anzustossen.

Facettierte Suche arbeitet nach dem Prinzip «Click & Filter», d.h. das Anklicken bestimmter Klassen in den angebotenen Facetten schränkt die Menge der infrage kommenden Ergebnisse nach und nach ein. **Bild 3** zeigt ein Beispiel aus dem EU-Verbundvorhaben Mace [4], an dem das DFKI mitgewirkt hat und in dem Information zu Architektur in ganz Europa über Facetten zugänglich gemacht wird. Mace bietet Facetten zu verschiedenen Repositorien, Sprachen, Ressourcen, Architekturkategorien und Kompetenzen an. Durch Anklicken einer Klasse, z.B. der Sprache «English» wird der Suchraum der anderen Facetten eingeschränkt und gleichzeitig nach Anzahl der verfügbaren Ressourcen sortiert. Die Auswahl des Repositoriums «Dynamo» führt zur weiteren Einschränkung der verbleibenden Optionen. Bei diesem Vorgehen erhält der Benutzer bei jedem Suchschritt einen exzellenten Überblick darüber, welche Eigenschaften die Ergebnismenge besitzt bzw. welche Eigenschaften zur Verfeinerung der Suche weiter zur Auswahl stehen.

### Visuelle semantische Suche

Die Visualisierung von Anfrageoptionen und Suchergebnissen ist durchaus sinnvoll, wenn man die inhärente Komplexität von Bedeutung betrachtet. Ausgewählte Visualisierungstechniken kombinieren daher unterschiedliche Darstellungsmetaphern wie Farbe, Grösse, Intensität, Abstand oder Richtung und helfen so, den Überblick über die Ergebnismenge nicht zu verlieren bzw. die Optionen für eine weitere Suche/Navigation zu erkennen. Auch an dieser Stelle kommen, ähnlich wie bei der facetierten Suche, Klassifikationssysteme zum Einsatz, um Wissen darzustellen und durch geeignete Interaktionsmöglichkeiten den Zugang dazu anzubieten. Bekannte Instanzen oder Dokumente werden je nach ihrer Relevanz in der Wissensbasis, der semantischen Nähe zur Anfrage, ihre Klassenzugehörigkeit oder eines betrachteten Beziehungstyps in einem Informationsraum dargestellt und können mit der Mouse-Over-Funktion weiter erschlossen werden.



