

12. GI-Fachtagung für Datenbanksysteme in
Business, Technologie und Web (BTW 2007)
5. bis 9. März 2007 - Aachen, Germany
<http://www.btw2007.de/>

YAWN: A Semantically Annotated Wikipedia XML Corpus

Ralf Schenkel Fabian Suchanek Gjergji Kasneci
Max-Planck-Institut für Informatik, Saarbrücken, Germany
{schenkel, suchanek, kasneci}@mpi-inf.mpg.de

Abstract: The paper presents YAWN, a system to convert the well-known and widely used Wikipedia collection into an XML corpus with semantically rich, self-explaining tags. We introduce algorithms to annotate pages and links with concepts from the WordNet thesaurus. This annotation process exploits categorical information in Wikipedia, which is a high-quality, manually assigned source of information, extracts additional information from lists, and utilizes the invocations of templates with named parameters. We give examples how such annotations can be exploited for high-precision queries.

1 Introduction

1.1 Motivation

Much of the existing work on structure-aware XML retrieval [AY⁺02, S⁺05, Sch02] has anticipated the existence of a huge number of heterogeneous XML documents with descriptive (i.e., self-explaining and semantically rich) tags. This has always been predicted as the natural consequence of the flexibility and variety of XML, where every author of a Web page can invent a different schema for her data on the Web. Given a semantically rich structural query like `//person[about(//work, physics) and about(//born, Germany)]`¹ (i.e., find people who work in physics and were born in Germany), such search engines would either consider structural similarity metrics of the query and documents [AY⁺02, Sch02] or semantic similarity of the tags used in the query and in the documents, using an ontology as background knowledge [S⁺05].

However, this revolution still has to happen. XML, even though widely used nowadays, usually is either generated from a structured database or used to store textual information with some structure. In the latter case, documents are content-rich, but tags are generic, while in the former case, there may be meaningful tag names, but there is often not much textual content. Therefore, today's typical XML search engines either ignore structural information completely, focussing on keyword queries alone, or deal with semantically weak structures as in `//article[about(., XML)]//section[about(./paragraph, retrieval)]`.

¹This query is formulated in NEXI, the query language of the INEX benchmark (<http://inex.is.informatik.uni-duisburg.de/2006/>), but could be expressed similarly in XPath with FullText extensions.

This paper bridges the gap between semantics-aware XML search engines and real world data by adding semantics to XML data, extending the ideas already used in Sphere-Search [G⁺05] and, more recently, by Chu-Carroll et al. [CC⁺06]. Based on the well-known and widely used Wikipedia collection, this paper presents *YAWN*², showing how to convert Wikipedia pages to XML that is annotated with concepts from WordNet, resulting in a huge annotated XML corpus with semantically rich tags. We exploit three sources of semantics: categorical information that has been added to most pages by the author (like “Albert Einstein is a physicist”), lists of similar pages, which are a common concept in Wikipedia (like lists of actors, songs, and companies), and invocations of predefined templates with human-readable parameters that encode factual information.

1.2 Related Work

There are a number of database-oriented XML benchmarks such as XMark [S⁺02] and XMach [BR01] that focus on performance aspects of XML databases. They usually provide no meaningful content and are therefore not suited for information retrieval. The INEX community has produced two XML benchmark collections, namely the IEEE collection and the Wikipedia collection [DG06], which combine huge corpora with queries, lists of relevant results, and a methodology to evaluate the quality of search results. However, these corpora do not include heterogeneous and semantically rich tags. The INEX Heterogeneous Track has started to collect different XML collections, but there are no self-explaining tags yet. A first step towards richly annotated data was made in the INEX Multimedia Track with the Lonely Planet collection³.

There have been several attempts to define a standard XML format for documents stored in Wikis, the most recent ones being the Wikipedia DTD⁴ that proposes a set of tags similar to HTML, and DocBook XML Export⁵ that defines XML tags for a subset of the Wiki markup. However, to the best of our knowledge, none of these covers the complete feature set of Wiki markup and is publicly available.

The Semantic Web community has recently launched a number of projects to add semantics to Wikis [Aum05, BG06, Sou05, V⁺06] which typically aim at adding semantic information at design time. In contrast, our approach exploits information that is already present in Wikipedia pages, without the need for any user interaction.

Information extraction from text and HTML data is an area with intensive work. Agichtein [Agi05] gives a survey of information extraction techniques with a focus on scalability in large collections with millions of documents. Approaches in the literature mostly either follow a rule-based paradigm [AGM03, CMM02, G⁺04], or employ learning techniques and/or linguistic methods [AFG03, CS04, Cun02]. Our algorithms, unsupervised and statistical in their nature, fit in the second class. The list extraction method used in Section 3.3 is similar to list extraction methods used in other information extraction sys-

²Yet Another Wikipedia Annotation project

³<http://inex.is.informatik.uni-duisburg.de/2005/tracks/media/index.html>

⁴http://meta.wikimedia.org/wiki/Wikipedia_DTD (last change APR-09-06)

⁵http://meta.wikimedia.org/wiki/DocBook_XML_export (last change JUN-19-06)

tems like KnowItAll [E⁺05] and SCRAP [FFT05]. Rule-based information Extraction from XML documents has been considered by Abolhassani et al. [AFG03]; in contrast, we follow a purely automatic approach. Annotated XML collections and their use for information retrieval were considered by Graupmann et al. [G⁺05] and Chu-Carroll et al. [CC⁺06]; the techniques presented in this paper perfectly integrate with these works.

2 Converting Wiki Markup to XML

2.1 Wiki Markup

With more than 1,4 million articles as of October 2006, Wikipedia⁶ is the largest general purpose encyclopedia that is freely available. The content of pages in Wikipedia, like in other Wikis, is formulated in *Wiki markup*, a combination of standard HTML tagging with specific constructs for structuring, tables, links etc. There is no formal specification of the language and its semantics yet. Figure 1 shows a simple example for a document in this language. Important building blocks of the Wiki markup language include structuring text by defining several levels of sections, different levels of emphasis for text parts, bulleted and numbered lists of different nesting depths, tables with rows, columns, headers and captions, links within the Wikipedia collection and to the Web, and inlined images.

```
==Introduction==
''Wiki markup'' is used in [[Wikipedia]].

==Language Components==
* tables
* lists
* and a lot more

==See also==
[http://www.wikipedia.org]
```

Figure 1: An example Wiki Markup page

In addition, Wiki markup can be arbitrarily mixed with HTML tags (which is most often used for layout purposes), and some Wiki markup symbols (like images and tables) may contain additional layout hints. Wiki markup also provides a template language where invocations of a template TEMP (in the form {TEMP}) are replaced with the definition of this template. Template invocations may also include values for parameters that then replace the parameter in the template's definition. A Wiki2HTML converter (written in PHP) generates HTML pages (that neither have semantically rich tags nor are guaranteed to be well-formed XML) from the Wiki markup input.

The different components of Wiki markup can be combined almost arbitrarily. However, this huge flexibility is at the same time a big problem, as authors often tend to overstrain the features. As an example, many authors make frequent use of tables for layout, resulting in

⁶<http://www.wikipedia.org/>

a deep nesting of tables and (sub-)sections. Secondly, the fault-tolerant converter does not encourage correct markup; while this is adequate to create HTML, this makes it difficult to create well-formed XML.

2.2 Generating XML

This section shows how we generate XML from the Wiki markup of the Wikipedia pages. We focus on the content of the pages, accepting that some layout information is lost in this process. The complete textual content of Wikipedia is available as a huge XML file (as of April 2006, this was about 6GB), which contains one element for each Wikipedia page with some meta information and its content in Wiki markup.

Our Wiki2XML converter runs in two phases, where each phase corresponds to a SAX-style scan of the input XML document. In the first phase, the existing pages are collected and redirections (i.e., pages that are simply links to other pages) are resolved. In the second phase, an XML document is generated for each page, which consists of three parts: (1) the preamble that sets the character encoding and includes a pointer to an XSLT for presenting the page, (2) the `article` element with its `header` child, which in turn has children that specify meta data like the page title and id, the last revision, and the categories of this page, and (3) the `body` child of the `article` element that contains the XML representation of the page's content. The conversion of the Wiki markup to XML is done as follows:

- Sections, subsections etc. are converted to `section`, `subsection` etc. elements, each with an `st` child that contains the section title.
- Both numbered and bulleted lists are converted to `list` elements with `entry` children corresponding to the different list entries. Nested lists are represented by `subentry`, `subsubentry` etc. tags.
- Tables are represented by a `table` tag. Inside this tag, there is one `row` element for each row of the table, which in turn has one `col` tag for each column. For header rows, the tag `header` is used.
- Links to other pages in Wikipedia are converted to `link` elements that correspond to an `XLink` to the target page's XML version. Links to web pages are converted to `weblink` tags with an `XLink` to the link target.
- Links to images are converted to `image` elements that correspond to an `XLink` to the image file located in a directory derived from the image name's MD5 hash.
- Markup for emphasis is converted to `b` and `it` elements.

Figure 2 shows the XML generated for the example Wiki markup from Figure 1. The exact conversion of Wiki markup to well-formed XML is more complex (and in fact sometimes impossible without human interaction) if the markup is mixed with arbitrary HTML tags. This is especially a problem with some tables that are defined in a mixture of HTML and

Wiki markup, but also simple tags like `
` render a generated document malformed. To solve this problem, we eliminate all HTML tags from the Wiki markup in a preprocessing step. As those tags are typically used only for layout purposes, we do not lose any semantics, but can generate much better XML. Additionally, as the original Wiki2HTML converter is quite tolerant to syntax errors, there is a substantial number of Wikipedia pages with syntactic problems⁷. Each document is tested for well-formedness after it was generated, making sure that the collection contains only syntactically correct documents.

```
<?xml version="1.0" encoding="UTF-8" ?>
<?xml-stylesheet type="application/xml" href="../../wikipedia.xslt"?>
<article xmlns:xlink="http://www.w3.org/1999/xlink/">
  <header>
    <title>Wiki markup</title>
    <id>42</id>
    <revision>
      <timestamp>2006-10-05 14:22</timestamp>
    </revision>
    <categories>
      <category>Markup languages</category>
    </categories>
  </header>
  <body>
    <section>
      <st>Introduction</st>
      <p><b>Wiki markup</b> is used in
      <link xlink:href="../../Wi/Wikipedia.xml" xlink:type="simple">
        Wikipedia
      </link>.</p>
    </section>
    <section>
      <st>Language Components</st>
      <list>
        <entry>tables</entry>
        <entry>lists</entry>
        <entry>and a lot more</entry>
      </list>
    </section>
    <section>
      <st>See also</st>
      <weblink xlink:href="http://www.wikipedia.org" xlink:type="simple">
        http://www.wikipedia.org
      </weblink>
    </section>
  </body>
</article>
```

Figure 2: Generated XML for the Wiki Markup from Figure 1

The XML dialect we use is similar to the one used by the INEX Wikipedia collection. However, we have better support for some details, including nesting of sections (in INEX, all nesting levels of sections are mapped to `sec` elements). Additionally, as the INEX

⁷The Wiki Syntax Project (http://en.wikipedia.org/wiki/Wikipedia:WikiProject_Wiki_Syntax) lists 8,400 for the April 21, 2005 dump, including 50 defective section headings, 80 HTML tags, 250 tables, 600 double quotes, 950 triple quotes, and 3400 square brackets.

Wikipedia collection has kept the HTML tags, it contains a noticeable amount of malformed or otherwise abnormal XML pages.

3 Semantic Annotation of Wikipedia Pages

We aim at finding high-quality semantic annotations for Wikipedia pages by a combination of exploiting manually created, high-quality information from categories that are assigned to most pages, and deriving additional information from highly structured documents such as lists (of persons, locations, etc.). To make the results applicable for a large suite of applications, we find annotations within the scope of a predefined ontology; we use WordNet [Fel98], the currently most extensive and widely used general-purpose thesaurus for the English language, but the results are transferable to any hierarchical ontology. The ontology provides us with a standard vocabulary for the annotations than can also be exploited for querying the annotated documents. Additionally, as our annotation algorithms are heuristic, we explicitly maintain an estimated confidence in the annotations.

3.1 Overview of WordNet

WordNet [Fel98] is a fairly comprehensive common-sense thesaurus carefully handcrafted by cognitive scientists. WordNet distinguishes between *words* as literally appearing in texts and the actual *word senses*, the concepts behind words. As of the current version 2.1, WordNet contains 81,426 synsets for 117,097 unique nouns. Often a single word has multiple senses, each of which comes with an estimation of its commonality and a brief description and is also characterized by a set of synonyms, words with the same sense, called *synsets* in WordNet. In this paper, we use the term *concept* for word senses, hence each concept corresponds to exactly one synset. WordNet provides relationships between concepts like hypernyms (i.e., broader senses), hyponyms (i.e., more narrow senses), and holonym (i.e., part of) relationships; for this paper, we focus on hypernym relationships.

Conceptually, the hypernym relationship in WordNet spans a directed acyclic graph with a single virtual source node 'ROOT' (that we introduced to get a connected graph) and seven first-level basic synsets (*entity*, *state*, *abstraction*, *event*, *act*, *group*, *possession*) that are children of the source node. Figure 3 shows an excerpt of that graph. For each concept in WordNet, there exists at least one, but usually several distinct root-to-concept paths (like for the concept 'singer' in the excerpt).

3.2 Exploiting Categories

The majority of Wikipedia pages is assigned to one or multiple categories. The page about Albert Einstein, for example, is in the categories `German_language_philosophers`, `Swiss_physicists`, and 34 more. Not all categories, however, imply that the entity

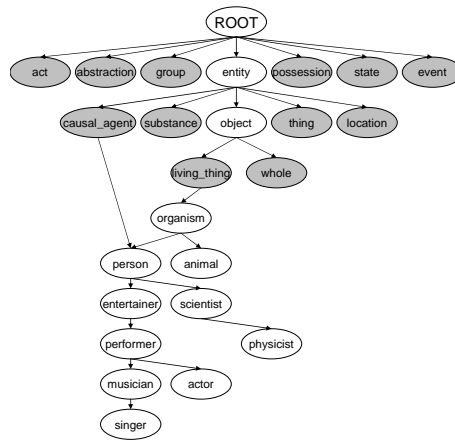


Figure 3: Excerpt of the WordNet DAG

described on the Wikipedia page is an instance of some concept. Some categories serve administrative purposes (like `Articles_with_unsourced_statements`), others yield non-conceptual information (like `1879_births`) and again others indicate merely thematic vicinity (like `Physics`). The administrative and non-conceptual categories are very few (less than a dozen) and can be excluded by hand. To distinguish the conceptual categories from the thematic ones, we employ a shallow linguistic parsing of the category name. For example, a name like `Naturalized_citizens_of_the_United_States` is broken into a pre-modifier (`Naturalized`), a head (`citizens`) and a post-modifier (`of_the_United_States`). Heuristically, we found that if the head of the category name is a plural word, the category is most likely a conceptual category. We used the Pling-Stemmer [S⁺06] to reliably identify and stem plural words. This gives us a (possibly empty) set of conceptual categories for each Wikipedia page.

As we want to use WordNet as foundation for our annotations, every conceptual Wikipedia category has to be linked to a corresponding WordNet concept. We experimented with different heuristics, including context-aware [S⁺05] and compactness-based [M⁺05] methods, and discovered that the simplest heuristics yields the correct link in the overwhelming majority of cases: We determine the WordNet concepts that the head of the category name refers to and link the category to the most common concept among them.

3.3 Exploiting Lists

Wikipedia contains many lists, which are an extensive, manually created and therefore high-quality source of information. In the Wikipedia Snapshot from April 2006 that we used for our experiments, there are 18,436 different lists. We consider the XML versions of the Wikipedia lists, i.e., the output of the XML generation process described in Section 2; this allows a better handling of structure in the lists. As an example, Figure 4 shows

an excerpt of the list of Germans from Wikipedia. To uniquely identify the elements of such a list, we assign a unique XPath expression to each element that consists only of name tests, child axes and position predicates. In the example, the `link` element that is a child of the second entry element in Figure 4; it is identified by the XPath expression `/article[1]/body[1]/section[1]/list[1]/entry[2]/link[1]`. Note that Wikipedia lists are not always designed as nested Wiki markup lists; there are many lists that are in fact tables. Our algorithm is not limited to Wiki markup lists and supports any type of regular structure.

It is evident that the example list is well structured, and it would be easy for a human to find out that all links point to pages about actors or, more generally, persons. A manual approach to exploit lists would require that a user identifies patterns in a list (either explicitly by an XPath expression or by highlighting them in an interface) and assigns them to a WordNet concept like 'actor'. However, while such a manual approach could be useful for small subsets of all lists, this tedious task does not scale to the whole Wikipedia collection.

```
<article>
...
<body>
  <section>
    <st>Actors</st>
    <list>
      <entry>
        <link xlink:href="..Ma/Mario+A$dorf.xml">
          Mario Adorf</link>, (born 1930), actor
        </entry>
      <entry>
        <link xlink:href="..Ha/Hans+A$lbers.xml">
          Hans Albers</link>, (1891-1960), actor
        </entry>
      <entry>
        <link xlink:href="..Mo/Moritz+B$leibtreu.xml">
          Moritz Bleibtreu</link>, (born 1971), actor
        </entry>
      ...
    </list>
  </section>
</body>
</article>
```

Figure 4: Excerpt from the XMLified Wikipedia list of Germans

3.3.1 Automatic Grouping of XPath Expressions

We propose an automated algorithm that exploits the fact that many pages have already been annotated with concepts derived from their categories, and only some pages are left to annotate. The algorithm, a variant of previously proposed algorithms for list extraction like [FFT05], proceeds in three steps: (1) it identifies parts of the list that are structurally similar (so-called *group candidates*), (2) it selects those group candidates where a large fraction of links points to pages with coherent annotations (so-called *groups*), and (3) it finds annotations that are common among them, and heuristically assigns these annotations to all pages in the group. In the example, if all but the third link point to pages that are labeled as 'actor', the algorithm would assign that concept to the page of 'Moritz Bleibtreu' as well.

In a preprocessing step, we first temporarily extend the annotations of all pages that are linked in the list. For each concept annotated to a page, we add all concepts on the root-to-concept paths of this concept in WordNet to the annotations of a page. This allows us to identify different annotations that are similar at a higher level of abstraction in WordNet, like 'actor' and 'singer', which are both subconcepts of 'performer'.

Group candidates are identified by grouping elements with similar XPath paths together. We say that two elements have a similar XPath if both paths have the same tag sequence and differ only in a single position. We label each group candidate with an XPath pattern that has the same tag sequence and the same positions as the elements in the group candidate, but a wildcard '*' at the position where the elements differ. As we maintain annotations only for pages (which are identified by link elements in the lists), it is sufficient to consider only group candidates where the last tag is link. We identify each XPath in a group candidate with the page its link element points to. In the example, the group candidate including the three link elements has the label `/article[1]/body[1]/section[1]/list[1]/entry[*]/link[1]`. We eliminate group candidates that are too small, i.e., consist of less than 5 elements.

To determine if a group candidate with n elements is a proper group, we count, for each concept c , the number f_c of times where c occurs as an annotation in the group. We accept the group if there is at least one concept where $\frac{f_c}{n} \geq \delta$, i.e., where the concept occurs in the annotations of at least a (configurable) fraction of all pages. Note that n includes pages without any annotations. Setting δ close to 1.0 gives a higher annotation quality, but may reduce recall, whereas setting δ close to 0 incurs a high danger of wrong annotations. In our experiments, $\delta = 0.75$ yielded good results. Each such concept c is then assigned to all pages in the group that are not already annotated with c ; the confidence of the annotation is set to $\frac{f_c}{n}$.

3.3.2 Outlier Detection

Even though most lists have a regular structure, sometimes outliers occur, i.e., small glitches in an otherwise perfectly regular structure. As an example, consider the excerpt of a list of songs shown in Figure 5. While, for most entries, the first link points to the singer's page and the second to the song's, this regularity is broken for one song that has two singers (shown in boldface). If we apply our algorithm in this setting, 'David Bowie' would be accidentally annotated as 'song'. Similar outliers can be caused by omitting links to pages that do not yet exist in Wikipedia (like some unknown singer) or by mistakes of the page editor. We apply a simple heuristic to detect such outliers that works as follows. First, we manually define the compatibility of a carefully selected set of base concepts towards the top of the WordNet DAG; these concepts are shown in grey in Figure 3. Compatible sets of base concepts are `{living,thing,causal_agent,group}` and `{whole,thing}`, all other combinations of base concepts are incompatible. Whenever a page p should be annotated with a new concept c , all base concepts $B(c)$ on paths from c toward the root node are computed and compared with the base concepts that have been assigned to p earlier in the process. The new concept is assigned to p if and only if each base concept in $B(c)$ is compatible with each already assigned base concept of p . This heuristic cleaning

eliminates most erroneous annotations.

```
<list>
  <entry>
    <link xlink:href="../../Jo/John+L$ennon.xml" xlink:type="simple">
      John Lennon</link> :
    <link xlink:href="../../Im/Imagine+(song).xml" xlink:type="simple">
      Imagine</link>
  </entry>
  <entry>
    <link xlink:href="../../Ne/Nena.xml" xlink:type="simple">
      Nena</link> :
    <link xlink:href="../../99/99+L$uftballons.xml" xlink:type="simple">
      99 Red Balloons</link>
  </entry>
  ...
  <entry>
    <link xlink:href="../../Qu/Queen+(band).xml" xlink:type="simple">
      Queen</link> &
    <link xlink:href="../../Da/David+B$owie.xml" xlink:type="simple">
      David Bowie</link> :
    <link xlink:href="../../Un/Under+P$ressure.xml" xlink:type="simple">
      Under Pressure</link>
  </entry>
</list>
```

Figure 5: Example for an outlier in a list of songs

3.4 Adding Semantic Tags to Pages

Annotations characterize a Wikipedia page and therefore should be stored with the page, enabling queries of the form “find pages of singers that...”. We therefore add tags that correspond to the annotations for a page right after the `article` element (see Figure 6). The tag names are derived from the WordNet synset that correspond to the annotated concept. The tags are augmented with the confidence, the ID of the WordNet concept, and the source of the annotation. In a NEXI-style query language, the example query fragment would be posed as `//singer[about(. . .)]`, exploiting the new annotation.

```
<article>
  <group confidence="1.0" wordnetid="26729" source="categories">
    <artist confidence="0.75" wordnetid="9187509" source="3 lists">
      <header>
        <title>Queen (band)</title>
        <id>42010</id>
        ...
      </header>
    </artist>
  </group>
</article>
```

Figure 6: Excerpt from the semantically annotated ‘Queen’ page

At the same time, the annotations of a page are also an important source of information in other pages that link to the page; this can be exploited for queries like “find concerts where the band Queen played”. To support this, we add the same tags also to links to the page in

other pages (see Figure 7). Once we have such an annotation of links, the example query fragment could be formulated as `//concert [about (//band, 'Queen')]`, exploiting the fact that any link to the Queen page will be annotated as `band` (as opposed to links to the Queen Mary ship or Queen Elizabeth II. of England).

```
<group confidence="1.0" wordnetid="26729" source="categories">
  <artist confidence="0.75" wordnetid="9187509" source="3 lists">
    <link xlink:href="../../Qu/Queen+(band).xml" xlink:type="simple">
      Queen</link>
    </artist>
  </group>
```

Figure 7: Semantically annotated version of a link to the 'Queen' page

4 Exploiting Implicit Semantics of Template Invocations

A rich source of semantics that is already included in Wikipedia are template invocations. However, unlike in the approaches presented in the previous sections, we exploit template invocations not for annotating a Wikipedia page as a whole, but for annotating pieces of information on that page. Templates are often used to generate a standard layout for structured information that is common to many pages. As an example, Figure 8 shows the invocation of the `Infobox_band` template that generates a table with some standard information on a musical band; this information is provided as parameters to the template. Similar templates exist for persons, countries, companies, rivers, software, and many more.

```
{{Infobox_band |
band_name      = Queen |
image          = [[Image:Queen.png|250px|right]] |
years_active   = 1971 - Present |
status         = Active |
country        = [[United Kingdom]]
}}
```

Figure 8: Example call of the `Infobox_band` template

For most templates, the name of a parameter is a clear indication of its semantics. We therefore exploit template invocations in a Wikipedia page to enrich the generated XML document with semantic annotations based on the template parameters. To do so, we try to map each parameter name to a WordNet concept, using the heuristics explained in Section 3.2. We then generate an element with the same name as the template. For each parameter, this element has a child element with the parameter name as name and the current parameter value as value. Figure 9 shows the XML that is generated for the example template invocation from Figure 8.

In contrast to the INEX Wikipedia collection that simply represents the template invocations and their parameters with a generic tag `<template>`, we believe that our approach is much more in the spirit of XML, allowing more natural XPath- and NEXI-style queries

```

<Infobox_band>
  <band_name>Queen</band_name>
  <image confidence="1.0" wordnetid="3782824" source="template">
    <imagelink xlink:type="simple"
      xlink:href="../../images/3/32/Queen.png"/>
  </image>
  <years_active>1971 - Present</years_active>
  <status confidence="1.0" wordnetid="13131686" source="template">
    Active
  </status>
  <country confidence="1.0" wordnetid="8023668" source="template">
    <link xlink:href="../../Un/United+K$ingdom.xml" xlink:type="simple">
      United Kingdom
    </link>
  </country>
</Infobox_band>

```

Figure 9: XML representation of the template call from Figure 8

such as `//article[about(., 'band')` and `contains(//country, 'USA')` and `contains(//status, 'active')`. Note that such queries can also be posed and answered if a user does not exactly know the schema, by relaxing tag names and other structural query conditions [AY⁺02, S⁺05, Sch02], or possibly by support of a DTD-aware graphical interface [vZ⁺06].

5 Applications

5.1 Concept-Based Information retrieval

An important application of annotations is concept-based retrieval. Graupmann et al. [G⁺05] have shown that annotating important classes of information like persons, locations, and dates can help to improve result quality, especially precision of results. As YAWN annotates with a huge set of diverse concepts from WordNet, it seems likely that these annotations can lead to further enhancements for the retrieval.

We have not yet done a thorough evaluation of the quality of annotations in general and their impact on result quality. However, to give a first impression of the use and effectiveness of annotations, we made some preliminary experiments with our XML search engine TopX [TSW05] on a small, annotated Wikipedia fragment, with NEXI-style [TS04] structured queries. We converted the first 10,000 Wikipedia documents (excluding redirections that contain only a pointer to another document) from the April 2, 2006 dump file into our XML format with semantic annotations, using the techniques presented in the sections before. The conversion failed for 49 documents, usually due to syntax problems in the input Wiki markup or unusual combinations of tables and sections. In our preliminary implementation, this took less than one hour on a standard notebook (the whole collection was converted in about 36 hours on a dual Xeon server machine).

We consider three types of queries: (1) queries that exploit only the annotation of complete

pages, optionally with content constraints, (2) queries that exploit the annotation of pages and links, optionally with content constraints, and (3) queries that additionally exploit annotations derived from template invocations. We collected 5 to 10 queries of each type and compared the performance of TopX with annotation-aware queries to plain keyword queries. Table 1 summarizes the results for the average precision@10 of the three query classes with and without annotation awareness. Note that we did not measure recall; we expect that some results will not be found due to errors in the annotation process. We now discuss some anecdotal results for the three query types.

Query type	Precision@10[annotations]	Precision@10[keywords]
1 (pages)	0.85	0.24
2 (pages+links)	0.96	0.14
3 (pages+links+templates)	1.0	0.0

Table 1: Experimental results with TopX

A simple example for the first query type are list-style queries like 'find scientists that won the Nobel prize' that would be formulated as keyword query `scientist nobel prize` without annotation awareness. If we exploit annotations, we can reformulate the query as `//scientist[about(., Nobel prize)]`, yielding a lot less nonrelevant results. For the example collection, TopX achieves a precision@10 of 0.2 for the keyword query, compared to 0.8 for the annotation-aware query. Here, the additional annotation serves as a prefilter for pages that simply mention the Nobel prize, but are not about scientists (like the page about Jimmy Carter who won the peace Nobel prize).

A typical example that shows the usefulness of the second query type is the query for musicians who have performed a song where 'space' occurs in the title, which could be formulated as `//musician[about(//song, space)]`. Without the annotation of the link, a search engine would find any occurrence of the term within the page, not only in the context of a song. As a result, the precision of the annotation query is perfect (1.0), whereas the corresponding keyword query finds no relevant results within the top 10.

The last query type is most powerful as it can exploit all three types of annotations. As an example, consider the query 'find mayors of towns in the (German state) Hesse'. While a keyword-based query has no chance to retrieve relevant results, the annotation-aware query `//town[about(//state, Hesse)]//mayor` that exploits tags introduced by the template `Infobox_Town_DE` yields only relevant results in the test collection. However, as not all city pages may use such a template, the recall is probably not perfect. To increase recall (at the cost of precision), using ontological tag expansion and structural similarity measures in the engine could help; this is subject of future work.

5.2 Other Applications

There are many more applications for an annotated XML corpus beyond the obvious information retrieval case sketched in the previous subsection, for example to help with

clustering and classification of Wikipedia pages. The annotations can further be exploited for query formulation, by letting a user pick concepts from WordNet and combine them to form a structured query, similar to the DTD-aware query formulation presented in [vZ⁺06]. The diverse structure introduced by the annotations will most certainly be a rich source for structural feedback [ST06] that exploits the structure of relevant results to generate more precise queries. And, finally, being the first real-life heterogeneous XML collection with both rich tags and content, we think that our new collection can serve as a large-scale benchmark for systems that exploit semantic annotation for retrieval, classification, clustering, etc, extending existing benchmarks with structural diversity.

6 Conclusions and Outlook

This paper presented YAWN, a project to create an XML version of Wikipedia with semantic information. We showed how to extract semantics from categories, lists, and template invocations, yielding a huge XML corpus annotated with semantically rich tags.

For future work, we plan to extensively evaluate the quality of the annotations and their effect for information retrieval, possibly within the INEX benchmark, and to consider some of the applications sketched in the previous section. Besides that, we will examine how we can integrate the collected information into our ontology and exploit it for other data sources, too. Finally, we plan to offer the annotated XML collection for public download.

References

- [AFG03] Mohammad Abolhassani, Norbert Fuhr, and Norbert Gövert. Information Extraction and Automatic Markup for XML Documents. In Blanken et al. [B⁺03], pages 159–174.
- [Agi05] Eugene Agichtein. Scaling Information Extraction to Large Document Collections. *IEEE Data Eng. Bull.*, 28(4):3–10, 2005.
- [AGM03] A. Arasu and Hector Garcia-Molina. Extracting Structured Data from Web Pages. In *SIGMOD 2003*, pages 337–348, 2003.
- [Aum05] D. Aumüller. Semantic Authoring and Retrieval in a Wiki. In *European Semantic Web Conference ESWC2005*, 2005.
- [AY⁺02] Sihem Amer-Yahia et al. Tree Pattern Relaxation. In *EDBT 2002*, pages 496–513, 2002.
- [B⁺03] Henk M. Blanken et al., editors. *Intelligent Search on XML Data*, volume 2818 of *Lecture Notes in Computer Science*. Springer, 2003.
- [BG06] M. Buffa and F. Gandon. SweetWiki: Semantic Web Enabled Technologies in Wiki. In *ACM 2006 International Symposium on Wikis*, pages 69–78, 2006.
- [BR01] T. Böhme and E. Rahm. XMach-1: A Benchmark for XML Data Management. In *BTW 2001*, pages 264–273, 2001.

- [CC⁺06] Jennifer Chu-Carroll et al. Semantic search via XML fragments: a high-precision approach to IR. In *SIGIR 2006*, pages 445–452, 2006.
- [CMM02] V. Crescenzi, G. Mecca, and P. Merialdo. RoadRunner: Automatic Data Extraction from Data-Intensive Web Sites. In *SIGMOD 2002*, page 624, 2002.
- [CS04] William W. Cohen and Sunita Sarawagi. Exploiting dictionaries in named entity extraction: combining semi-Markov extraction processes and data integration methods. In *KDD 2004*, pages 89–98, 2004.
- [Cun02] H. Cunningham. GATE, a General Architecture for Text Engineering. *Comput. Humanit.*, 36:223–254, 2002.
- [DG06] Ludovic Denoyer and Patrick Gallinari. The Wikipedia XML Corpus. *SIGIR Forum*, 40(1):64–69, June 2006.
- [E⁺05] Oren Etzioni et al. Unsupervised named-entity extraction from the Web: An experimental study. *Artif. Intell.*, 165(1):91–134, 2005.
- [Fel98] C. Fellbaum, editor. *WordNet: An Electronic Lexical Database*. MIT Press, 1998.
- [FFT05] Bettina Fazzinga, Sergio Flesca, and Andrea Tagarelli. Learning Robust Web Wrappers. In *DEXA 2005*, pages 736–745, 2005.
- [G⁺04] Georg Gottlob et al. The Lixto Data Extraction Project – Back and Forth between Theory and Practice. In *PODS 2004*, pages 1–12, 2004.
- [G⁺05] Jens Graupmann et al. The SphereSearch Engine for Unified Ranked Retrieval of Heterogeneous XML and Web Documents. In *VLDB 2005*, pages 529–540, 2005.
- [M⁺05] Dimitrios Mavroeidis et al. Word Sense Disambiguation for Exploiting Hierarchical Thesauri in Text Classification. In *PKDD 2005*, pages 181–192, 2005.
- [S⁺02] A. Schmidt et al. XMark: A Benchmark for XML Data Management. In *VLDB 2002*, pages 974–985, 2002.
- [S⁺05] Ralf Schenkel et al. Semantic Similarity Search on Semistructured Data with the XXL Search Engine. *Information Retrieval*, 8(4):521–545, December 2005.
- [S⁺06] Fabian M. Suchanek et al. LEILA: Learning to Extract Information by Linguistic Analysis. In *2nd Workshop on Ontology Population (OLP2) at ACL/COLING*, 2006.
- [Sch02] Torsten Schlieder. Schema-Driven Evaluation of Approximate Tree-Pattern Queries. In *EDBT 2002*, pages 514–532, 2002.
- [Sou05] A. Souzis. Building a Semantic Wiki. *IEEE Intelligent Systems*, 20:87–91, 2005.
- [ST06] Ralf Schenkel and Martin Theobald. Structural Feedback for Keyword-Based XML Retrieval. In *ECIR 2006*, pages 326–337, 2006.
- [TS04] Andrew Trotman and Börkur Sigurbjörnsson. Narrowed Extended XPath I (NEXI). In *INEX Workshop 2004*, pages 16–40, 2004.
- [TSW05] Martin Theobald, Ralf Schenkel, and Gerhard Weikum. An Efficient and Versatile Query Engine for TopX Search. In *VLDB 2005*, pages 625–636, 2005.
- [V⁺06] Max Völkel et al. Semantic Wikipedia. In *WWW*, pages 585–594, 2006.
- [vZ⁺06] Roelof van Zwol et al. Bricks: The Building Blocks to Tackle Query Formulation in Structured Document Retrieval. In *ECIR 2006*, pages 314–325, 2006.