

# Link Annotation

(Part of MRef Project)

Maciej Janik, Ravi Pavagada, Samir Tartir, Bilal Gonen  
Department of Computer Science,  
University of Georgia,  
415 Boyd Graduate Studies Research Center,  
Athens, GA 30602-7404

{mjanik, ravipr, starter, gonen}@uga.edu

## Abstract

Web pages in the web represent certain concepts in the domain they fall in, and the connections between them represent the relations between the concepts they represent. In the current web, people are using links blindly without knowing what these links point to, or what kind of relationship this link represents. With the advent of the semantic web, concepts and relationships among them are represented in an ontology. This can be utilized to make links more meaningful. Web pages can be searched, browsed or even reorganized based on their concept and relationship labels. Links in a webpage can render useful information about the page it is pointing to. We can annotate webpage & its links with appropriate concepts from ontology. This paper presents new idea of propagating concept from a webpage to the links pointing to that page or from the links to the webpage. Propagation of concepts is based on certain criteria which will be discussed later in this paper. We also propose a new idea of automated voting which is used to choose the right concept or relation from a number of concepts & relation matches.

## 1. Introduction

The network of hyperlinked documents, as it exists now, lacks semantic information in machine understandable form. It can only be browsed or searched by keywords -not concepts. There exist projects that automatically or semi-automatically annotate web pages with concepts taken from ontology. This effort makes web pages more understandable for machine processing and searching. In our project we would like to focus more on navigational implications of adding semantic annotation to web pages.

Currently user or machine navigates between web pages by traversing them via hyperlinks. Decision if accessed page is relevant to the undertaken search can be made only after retrieving and analyzing the destination web page. In our project we would like to add more semantic meaning to links themselves on the source page, so concepts included on target page can be evaluated without retrieving page itself.

In this paper, we use well formed computer science department ontology to annotate links and web pages with concepts. Web pages and links of the page can then be associated with concepts and relations from ontology. For example, web pages from computer science department of University of Georgia web site can be associated with concepts such as faculty, department, Course, lecturer, research assistant etc... These web pages can therefore be treated as concept instances. Relationship can be defined between a webpage and its link. For ex. A given student webpage might be having a link to his course page. In ontology there could be a relation say *'takes'* between the student & the course. This information will be annotated in the links along with the link concept *'course'*. We have used ontology dictionary which associates labels to each concept and relations in the ontology. These labels are very useful in concept matching. Labels form the key role, since they are matched with the page contents & link window to extract appropriate concepts. We haven't used NLP techniques to get the concept matches.

Our goal is to start with set of plain, connected web pages and by extracting information and matching them with the ontological concepts and also annotate the links with concepts and relations joining them. In this project we would like to utilize already known algorithms and solution for page annotations. We think that combining different approaches of page annotation and information/concept propagation between web pages can improve the overall quality of annotated data.

## Paper Outline

Paper is organized as follows. In Section 2, we describe the related work and our ideas. Section 3, briefs our work and discusses the architectural of the proposed system. Section 4, approach we took in building the proposed system. Section 5, describes propagation and voting schemes used. Section 6, describes the testing and experimental results, Section 7, discussion and analysis and Section 8, Conclusion and Future work.

## 2. Related Work & Our Ideas

There are papers on HTML Tag tree extraction or deriving link context, one of them is "Deriving link-context from HTML tag tree" by Gautam Pant et al. and other papers on automated semantic annotations. "SemTag & Seeker: Bootstrapping the semantic web via automated semantic annotations" talks about automation of web page annotations. "Mining the link structure of semantic web", by Souman Chakrabarti et al., talks about HITS algorithm which takes advantage of the hubs in some fields & uses techniques that take advantage of social organizations of the web and allocates weights for the hub pages & authorities in iterative process. The paper "On extracting link information by relationships instances from a website" by Myo-Myo Naing et al, talks about a web page which is being associated with a concept in ontology & links two different web pages based on the relationship between concepts in the ontology.

Our work is slightly different from their work. We incorporate voting of relations whenever there are more than one relation matches between two concepts. Concept matching is an area in itself & there are lots of papers on it. There are lots of AI & natural language processing techniques used to achieve this. As mentioned previously we have concentrated more on concept labels defined in the ontology to find a concept match. Our work concentrates more on the information which is around the link, i.e. link context & match the link to a concept. New idea of concept propagation is proposed which would propagate concepts from a Webpage to the links pointing to that page if there is a tie in the number of concept matches for the given set of links. Propagation from links to page is done if most of the links agree on a single concept. Voting of concepts & relations is done whenever there is ambiguity.

### 3. Architecture Overview at a high level

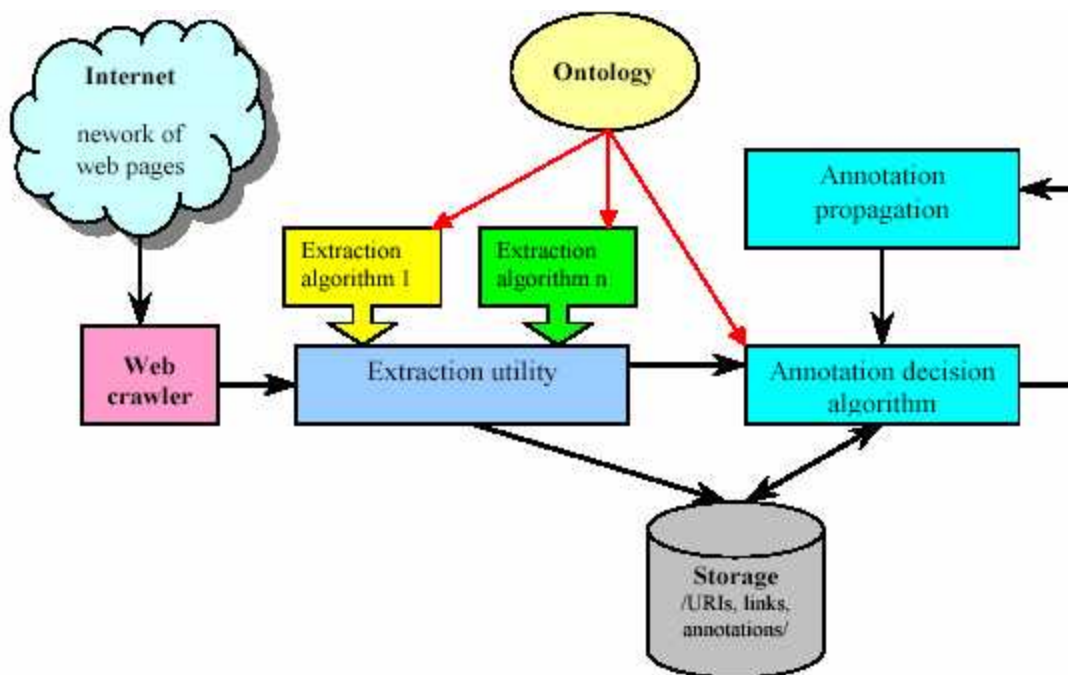


Figure 2 General system architecture

We would like to make our system modular and expandable for future needs. As we cannot modify the content of web pages, we can only keep discovered annotations of pages and links in snapshot of selected web pages.

#### WebCrawler

Web Crawler, crawls the web structure and supplies the raw data for further analysis. HTML from web pages is analyzed by extraction utility. Here extraction mechanism tries to match the whole page to some concepts in ontology.

## Ontology Dictionary

Ontology dictionary is the key part of our project. Dictionary labels are assigned for concepts and relations in the ontology. Dictionary labels for relations & concepts are comprised of hypernyms, synonyms & homonyms. We have added RA as a label for Research Assistant, TA as a label for Teaching Assistant etc.

```
<rdfs:Class rdfs:about="http://protege.stanford.edu/kb#AssistantProfessor"
rdfs:label="Assistant Professor">
<rdfs:label>Assistant Prof</rdfs:label>
<rdfs:label>Assistant Faculty</rdfs:label>
<rdfs:subClassOf rdfs:resource="http://protege.stanford.edu/kb#professor" />
</rdfs:Class>
```

We add labels to each of the concept in the ontology. Above is an extract from the ontology which describes an `rdfs:class` Assistant Professor and its associated labels namely Assistant Professor, Assistant Faculty etc.

## Extraction utility

Extraction utility is comprised of page & link analyzer, which analyses the page for the tags and assembles a vector of number of concept matches for each tag. The vector size is determined by the number of concepts in the ontology. We have prioritized various html tags in the webpage based on its importance. For example. `<Title>`, `<Head>`, and `<Body>` tag are given most importance. It also tries to categorize links in this web page based only on information contained in the link window. Link extractor extracts the text of information based on the window sizes or number of bytes of text before and after the link. It then assembles the concept weights based for each of the link window sizes namely 0, 50 (25 words above and below the link) and 150 (75 words above and below the link).

## Annotation decision

Voting is done whenever there is more than one concept matching a given page or a link. Based on the values set in the configurator, i.e. relative importance of tags or it could be based on relative size of the windows 0, 50 or 150 the voter calculates the new vote by calculating the product of the weight vector to the weights assigned in the configurator. Our configurator is flexible & easier to change. We have assigned weights of 0.5, 0.3, & 0.2 weights for anchor text window sizes of 0, 50 & 150 respectively. We have assigned weights of 0.4, 0.3, and 0.3 for "title", "head", and body" tags respectively

## Database Storage for persistence

Once the voting is done for the web pages and links we update the webpage and links table in the database. All extracted information is stored in persistent storage along with the matched concept for web page and its links. One advantage of our approach is that we have designed the project in such a manner that the all the web page and web link information are stored in the tables of our database. We crawl the web pages and load the tables in the database with the concepts. Then we follow the links to store its content and the relevant concept matches. After propagation tables of the web pages and the links are updated along with the concept and its relation matches.

## Propagation of concepts

As last, came the decision and propagation loop. Now the extracted information is analyzed again. Some of extracted information may be deleted from page; some can be inferred or pushed from links to page. In this step web pages are analyzed in network and we allow annotation flow between nodes. Both from page to describing link and from link to described page. This is iterative process and few iteration the network should reach some stable (or near-to-stable) state. In such state we would say that the selected network is annotated and can be used in semantic navigation. Propagation of concepts and voting is discussed will be discussed in detail in the coming sections.

## 4. Approach

The approach will work in two general phases: Preparation and Annotation.

### I. Preparation:

Here a deep analysis of the Computer Science department in the University of Georgia will be conducted, resulting in building an ontology that represents the current structure of the department.

This resulting ontology will be used in the next phase for annotation.

### II. Annotation:

This is where the actual process of page and link annotation will take place. This phase is divided into three stages:

1. Page annotation.
2. Link annotation.
3. Relationship annotation.

#### 1. Page annotation

In this stage, all the pages in the Computer Science department site will be analyzed in one of the current methods, or a new method that we might need to develop. The result of this analysis will be a mapping between a certain page, and a node in the ontology designed in phase I.

## 2. Link annotation

Here, each page will be scanned for links that point to pages in the same domain, and each link will carry the annotation of the page it points to.

## 3. Relationship annotation

This is the final stage that defines which type of relationship the link defines. This relationship is obtained from the ontology based on the types (concepts) of the page with the link, and page the link points to.

The resulted annotated pages will be stored in a database the application has access to write to and issue queries against.

# 5. Voting and Propagation

## Voting Algorithm

### I. Voting for Web pages

1. for each tag entry in the configurator do
2. begin
3. for each concept in the ontology ( since vector returned by StructuralNLP which is of the size of no. of concepts in the ontology)
4. begin
5. Calculate the weights based on the number of matches for each concept and weight assigned to the tag as a whole. This is important since the relevance of each tag is different. This was discussed earlier in the paper.
6. end
7. end
8. Select the maximum concept weight among all the vectors

### II. Voting for Web link

1. For each of the size of the hypertext window in the configurator do (size of the hypertext window described in the configurator are 0,50 & 150 and the weights associated with each of them are 0.5, 0.3, & 0.2 respectively.
2. begin
3. for each concept in the ontology ( since vector returned by StructuralNLP which is of the size of no. of concepts in the ontology)
4. begin
5. Calculate the weights based on the number of matches for each concept and weight assigned to each of the individual anchor text sizes. This is important since the relevance of each of the sizes of the anchor text window is different. This was discussed earlier in the paper.
6. end
7. end
8. Select the maximum concept weight among all the vectors

### III. Voting for Relations

1. Given two concepts
2. We get all the relations between the two concepts
3. Traverse all the concept nodes that are above a given two concept in the ontology & extract the relations between them
4. Then we match these relations obtained with the text surrounding the hypertext window.
5. We choose the concept that has the maximum number of keyword matches among all the hypertext window vectors.

The main drawback at this point with respect to relation voting is that we haven't concentrated on weighting relations between the concepts with different weights.

### Propagation Algorithm

Propagation of concepts is done to increase the accuracy in the concept matches for a given web page and its links. Propagation is done after the voting stage. At the end of the voting we would have the concepts for web pages and its links stored in the database.

Propagation of concept based on the following criteria:

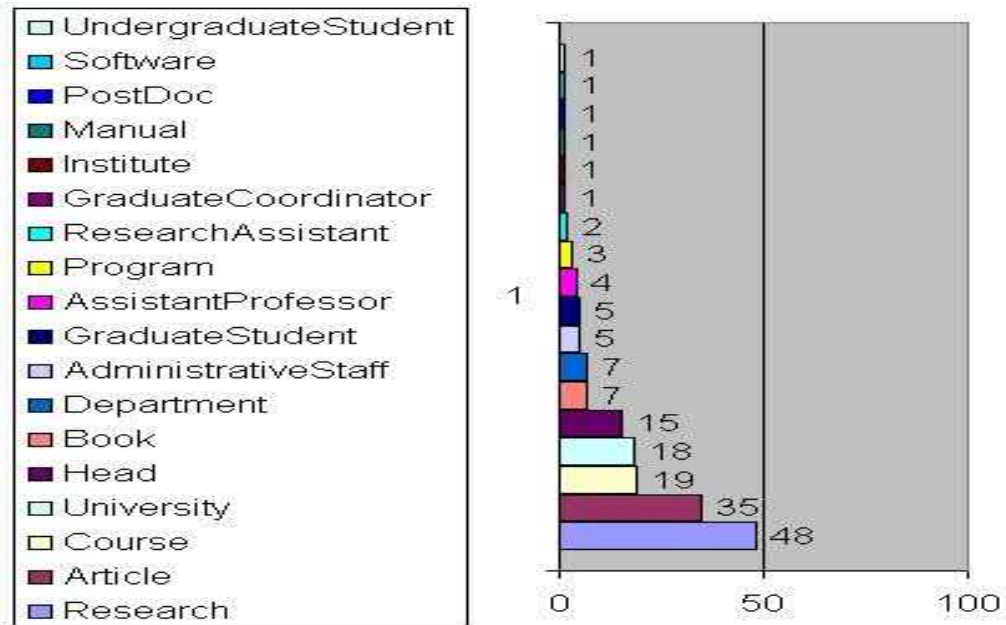
1. First, we get the concept match to a page querying the webpage table from database.
2. We then get the concept matches for the links pointing to a given page. (This is based on the link context). This is again got by querying the web link table in the database.
3. We propagate the concept from the links to page if there are maximum numbers of links matching to a given concept.
4. If there is a tie then we propagate the concept from the webpage to the links. (Reason behind this is based on the intuition that webpage concept has higher priority than the concepts of the tied links.)

### 6. Testing and Experimental results

Testing was done on [www.cs.uga.edu](http://www.cs.uga.edu) domain. Web crawler started crawling from <http://lsdis.cs.uga.edu> and was allowed to crawl only the CS domain. We ran our test crawling 200 web pages and the total number weblinks crawled were about 4599.

We used the following test cases:

1. Finding concept for page
2. Finding concept for link
3. Propagating down concept, i.e. from a webpage to the the links pointing to it.
4. Propagating concept up, i.e. propagation occurs from the weblinks to the webpage.
5. Finding relation for a given link



Above fig. describes the page concept and number of matching pages

Total concept matches for the pages crawled were about 174. Concepts weren't matched to a page since there wasn't any concept label matching to the pages.

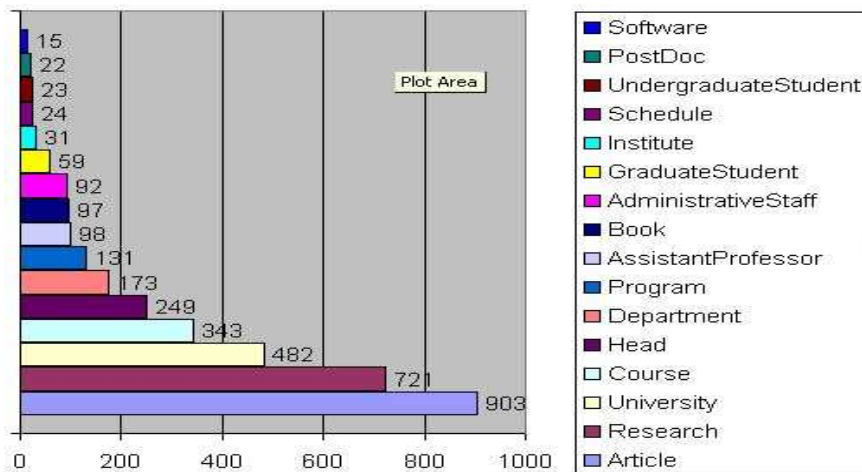


Fig. Describes the concept and number of matching links

There were about 3463 links matching to a concept out of a total of 4599.

Example for the best concept match

URL: <<http://www.cs.uga.edu/academics/UGProgram/4900courses.htm>>



Matching Concept: Course

Some of the good examples for concept matching before propagation

- <<http://webster.cs.uga.edu/~budak/>>  
concept: AssistantProfessor
- <http://webster.cs.uga.edu/~Ekochut/Research/>  
concept: Research
- <<http://lsdis.cs.uga.edu/Projects/METEOR-S/Downloads/>>  
concept: Article
- <<http://www.cs.uga.edu>>  
concept: Department

Propagation of Concepts

Some of these webpages initially had '*Research*' as their page concept. After propagating from links to the page we had the following results. Below are some of the examples of propagating concepts from links to the pages.

- Link concept (#Department) from <<http://www.cs.uga.edu/~jam>>
- Link concept (#Department) from <<http://lsdis.cs.uga.edu/~devp>>
- Link concept (#GraduateStudent) from <<http://lsdis.cs.uga.edu/~mperry>>
- Link concept (#GraduateStudent) from <<http://lsdis.cs.uga.edu/~aleman/>>
- Link concept (#Department) from <<http://lsdis.cs.uga.edu/~cthomas>>

- Link concept (#Department) from <<http://lstdis.cs.uga.edu/~kunal>>
- Link concept (#Department) from <<http://lstdis.cs.uga.edu/~kaarthik>>
- Link concept (#TeachingAssistant) from <<http://lstdis.cs.uga.edu/~mperry>>
- Concept matched to a page was #Research, and the propagated concept was #Department

#### For Relations

- The page <<http://lstdis.cs.uga.edu/about/index.php?page=1>> matched to a concept [#AdministrativeStaff] and <<http://www.uga.edu/>> matched to concept [#University]
- There was only one relation between them in the ontology. The relation found between these two concepts were <#works>

#### Summary of our experimental test cases:

In our experiment, first of all, we crawled 200 web pages. We set starting pages as <http://lstdis.cs.uga.edu> and <http://www.cs.uga.edu>. We limited the crawling area as in the "cs.uga.edu" domain. After crawling 200 pages, we crawled the links within each pages. Thus, the total number of link crawled were 4599. By running our algorithm on these 200 webpages, our algorithm assigned some concepts to 174 pages of them. Also running our algorithm on the 4599 weblinks, our algorithm was able to assign concepts to 3463 links.

Following are some of the examples,

< <http://webster.cs.uga.edu/~budak/> > concept matched: AssistantProfessor  
< <http://webster.cs.uga.edu/~Ekochut/Research/> > concept matched: Research  
< <http://lstdis.cs.uga.edu/Projects/METEOR-S/Downloads/> >  
Concept matched: Article  
< <http://www.cs.uga.edu> > concept: Department

Propagation algorithm used showed good results. Below is an example of its effectiveness.

The Computer Science web page had concept #Research before propagation. By applying the propagation method considering 93 pointing links, original page concept "Research" has been changed as concept "Department".  
Concept was propagated from Links to Page.

## 7. Conclusion and Future Work

Our approach of using ontology labels for relations & concepts in ontology was very beneficial in concept matching. We were able to match most of the web pages to the concept in the ontology. Labels were represented by hypernyms, synonyms & homonyms. Our propagation algorithm showed excellent results. We were able to compare the effectiveness of the algorithm by comparing it with concept before propagation. Voting was done based on the importance of individual tags and also based on the importance of the various anchor text window sizes. Relation voting seemed to work pretty well. Relation voting was done whenever there were more than one relation matches between two concepts. Our algorithm or methodology could be changed by adding different weights to relations between concepts, i.e. we traverse the ontology tree to find all the possible relations between the concepts by traversing the tree in a bottom up fashion. Since our algorithm uses these relations and matches the keywords around the anchor text window. Our future work would be to include different weights to the relations as we traverse the tree in a bottom up fashion. We still need to tune the ontology as there are no concept matches for some of the web pages crawled. Using various label names to a given concept may not be a very best idea compared to NLP techniques. Our ontology is not populated with instances so that we could use it for semantic web search. Initial experimental results were very promising, and we wish to work on this a little more.

## References

1. Effects of Link Annotations on Search Performance in Layered and unlayered Hierarchically Organized Information Spaces  
<http://lhncbc.nlm.nih.gov/lhc/docs/published/2001/pub2001040.pdf>
2. Mining the Link Structure of the WWW  
<http://citeseer.ist.psu.edu/chakrabarti99mining.html>
3. Deriving Link-context from HTML tag tree  
[http://dollar.biz.uiowa.edu/~pant/Papers/tagTree\\_dmkd.pdf](http://dollar.biz.uiowa.edu/~pant/Papers/tagTree_dmkd.pdf)
4. Search Engine-Crawler Symbiosis: Adapting to Community interest  
<http://dollar.biz.uiowa.edu/~pant/Papers/se-crawler.pdf>
5. Automatic resource compilation by analyzing hyperlink structure and associated text  
<http://marco.uminho.pt/disciplinas/UCAN/BD/Artigos%20Recomendados/chakrabarti98automatic.pdf>
6. On Extracting Link Information of Relationship Instances from a Web Site  
<http://citeseer.ist.psu.edu/662560.html>
7. Semantic Blogging and Bibliography Management  
[http://www.w3.org/2001/sw/Europe/reports/open\\_demonstrators/hp-requirements-specification.html](http://www.w3.org/2001/sw/Europe/reports/open_demonstrators/hp-requirements-specification.html)
8. Ontology-based Web Annotation Framework for Hyperlink Structures
9. A Study of User Model Based Link Annotation in Educational Hypermedia  
<http://www2.sis.pitt.edu/~peterb/papers/JUCS98.pdf>
10. Semantic Linking - A context based approach to Interactivity in Hypermedia
11. Web Document Searching Using Enhanced Hyperlink Semantics Based on XML  
<http://www.db-net.aueb.gr/hercules/papers/ideas.pdf>
12. Automatic Annotation of Content-Rich HTML Documents: Structural and Semantic Analysis <http://citeseer.ist.psu.edu/668883.html>
13. Semtag and seeker: bootstrapping the semantic web via automated semantic annotation  
<http://www.almaden.ibm.com/webfountain/resources/semtag.pdf>
14. Media-independent correlation of Information: What? How? -- MREF paper  
<http://www.computer.org/conferences/meta96/sheth/>