# CS 221 Project 3 M3 - Search Engine

Jun Ma 29846938
Chang Liu 47927446
Lu Chen 35888463

## 1.Improvements made for Search Engine

In the milestone 2, we realized the basic functions of our search engine. In this final stage of the Search Engine project, we studied different strategies to optimize the ranking of the search results as shown below:

- Generate new Inverted Index of anchor text in both one-gram and two-gram.
- Generate new 2-gram Inverted Index of both title and body of the document.
- Optimize the ranking quality based on weigh adjustment on inverted index of document & anchor text and Page Rank.
- Automation on NDCG computation and Google search result extraction.

Based on the previous two steps and improvements above, we get the NDCG sores before and after performance improvements as shown in table 1. Obviously, our search engine results are improved in almost all queries.

Table 1. NDCG Score before and after performance improvements

| Queries | NDCG before improvement | NDCG after improvement |
|---|---|---|
| mondego | 0.61 | 0.61 |
| machine learning | 0.26 | 0.33 |
| software engineering | 0.16 | 0.17 |
| security | 0.23 | 0.19 |
| student affairs | 0.24 | 0.41 |
| graduate courses | 0 | 0.48 |
| Crista Lopes | 0 | 0.41 |
| REST | 0.44 | 0.44 |
| computer games | 0 | 0.28 |
| information retrieval | 0.24 | 0.58 |

## 2: Ranking Algorithm

- Read in the query, apply two-gram algorithm or one-gram algorithm based on the query size.
- Match the query with the token in the inverted index, compute corresponding joint relevance scores from features of Anchor Text, TF-IDF of Title & Body text and PageRank respectively.
- In two-gram algorithm, List of matched document is divided into two buckets with the prior one containing the two-gram tokens and the other one containing either of the query word.
- Ranking inside the prior bucket are based on both two-gram and one-gram relevance scores from features of Anchor Text, TF-IDF of Title & Body text. Assigned these two parts with different weights and join them together with Page Rank values of the document set to finalized the total score.
- Ranking of the second bucket is same as the one-gram algorithm.
- Return top 10 results with highest relevance scores.

## 3: NDCG computation

We use Discounted cumulative gain (DCG) to measure the ranking quality of our search engine, the formula is shown below:

$$DCG_i = rel_i + \sum_{i=2}^{p} \frac{rel_i}{\log_2(i+1)}$$

DCG score is normalized taking Google search result as the ground truth. The graded relevance scale of top 5 documents in the search result set ranges from 5 to 1.

For example, when we search query "Crista Lopes" in Google, the top five results are :

*http://www.ics.uci.edu/~lopes*
*http://www.ics.uci.edu/about/search/search_dept_in4matx.php*
*http://www.ics.uci.edu/~lopes/patents.html*
*http://www.ics.uci.edu/faculty/profiles/view_faculty.php?ucinetid=lopes*
*http://mondego.ics.uci.edu*

When we search the same query in our search engine, the top five results are:

*http://www.ics.uci.edu/~lopes*
*http://www.ics.uci.edu/~lopes/patents.html*
*http://www.ics.uci.edu/~lopes/aop/aop-pics.html*
*http://www.ics.uci.edu/~lopes/aop/aop.html*
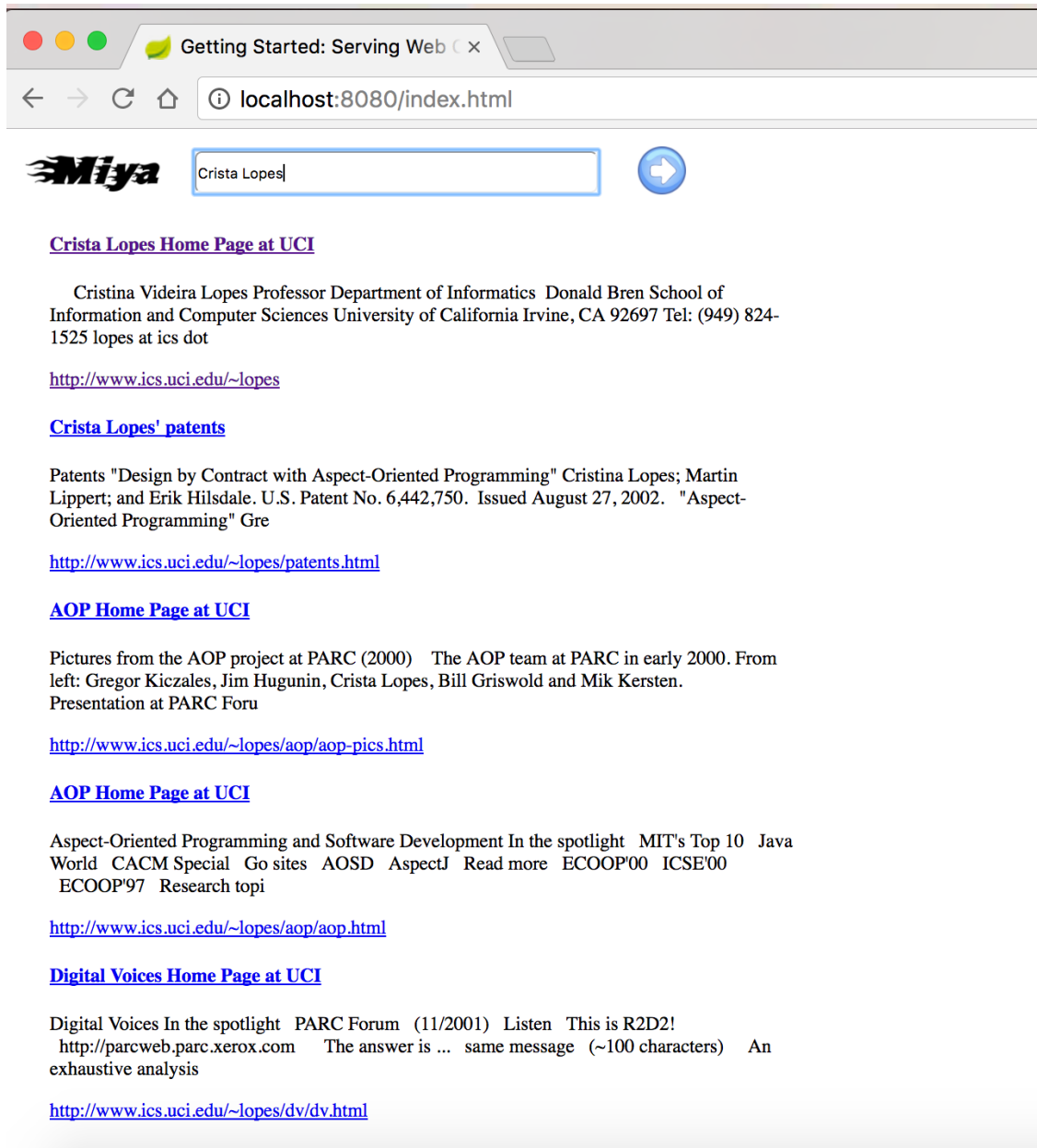*http://www.ics.uci.edu/~lopes/dv/dv.html*

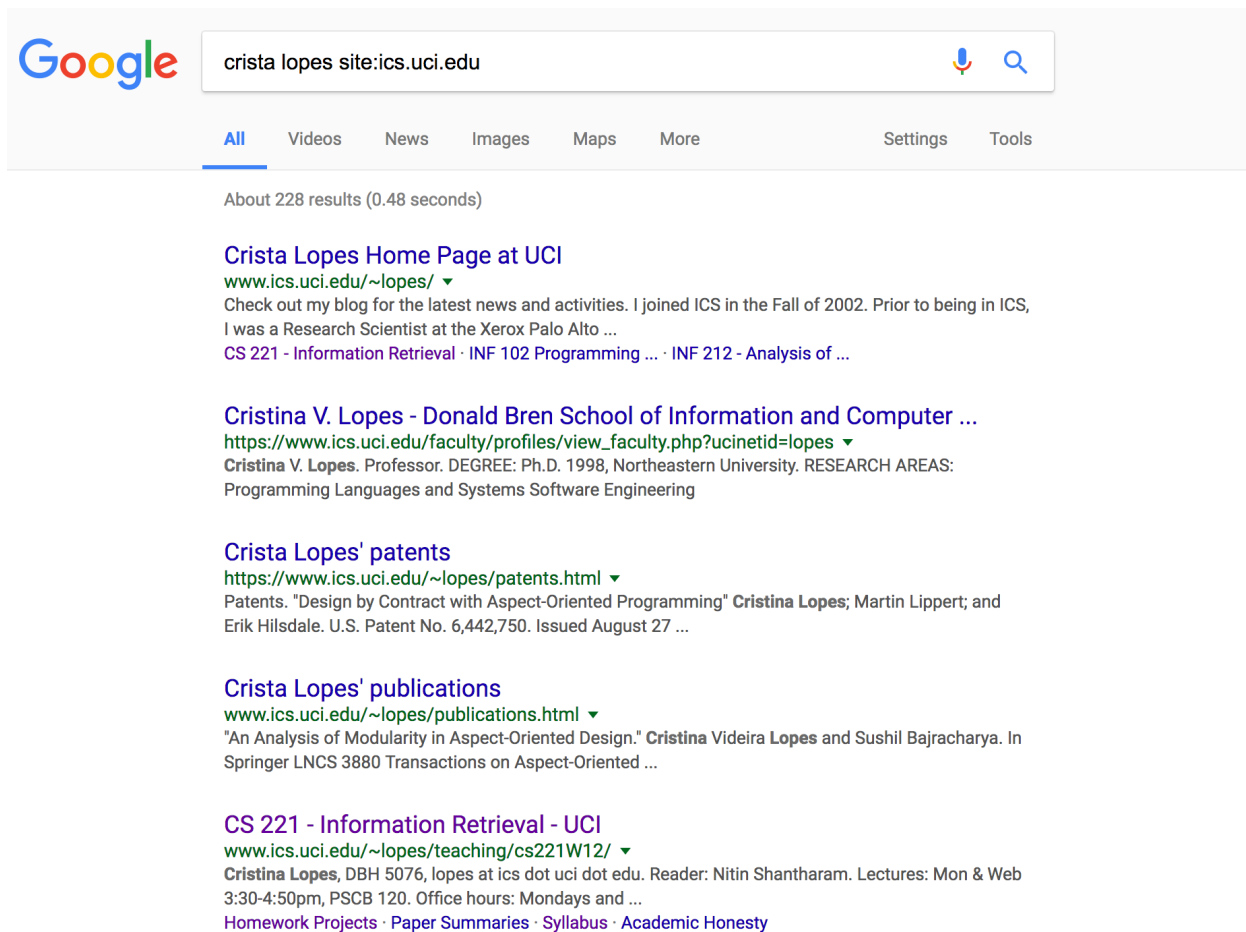Then we  The snapshot of NDCG result on query "Crista Lopes" is shown in picture below.

```
local DCG5:7.523719014285829
ideal DCG5:12.323465818787765
NDCG5:0.6105197291832893
```

# 4.Search Engine Result

The search result for query "Crista Lopes" from our search engine is shown below.



The corresponding search result from google is shown below.

## 5.Further Improvement

Based on all three steps of Project 3, we realized the search engine with relatively satisfying search results on given queries. For further improvement, the URL pool may be extended thus to cover the entire subdomain as needed. Also, to strengthen the PageRank factor by inserting more hub links can also help with precision. What's more, more file types should be supported, for instance, files in pfd or ppt play an important role in academic search.

In conclusion, we successfully implemented a complete search engine for the ICS domain. Specifically, we generated inverted index that maps words to documents, then constructed an interface for search queries and returned top 10 results with highest relevance score. Finally, optimized relevance algorithm is implemented and the ranking performance has been improved.