

# Ling 573 Summarization Presentation

**Thomas Marsh**  
University of Washington  
ling 573, Spring 2015  
sugarork@uw.edu

**Michael Roylance**  
University of Washington  
ling 573, Spring 2015  
roylance@uw.edu

**Brandon Gahler**  
University of Washington  
ling 573, Spring 2015  
bjg6@uw.edu

## Abstract

TODO: fill out

## 1 Introduction

TODO: fill out

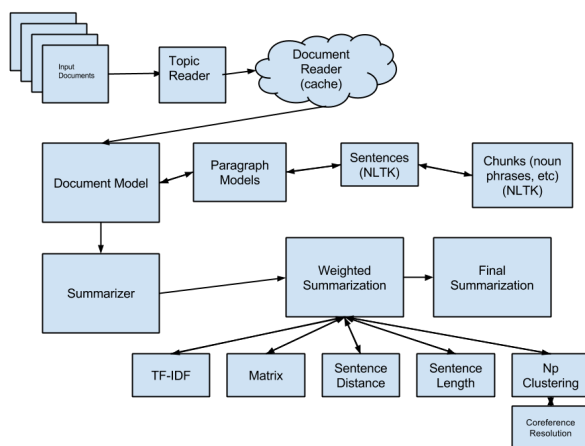
## 2 System Overview Diagram

TODO: fill out

## 3 Approach

### 3.1 System Architecture

Figure 1: System Architecture



### 3.2 Content Selection

We implemented a variety of methods initially to help with content selection. Each of these five separately implemented content selection systems produces a score for each sentence from 0.0 to 1.0. Empirically determined system weights are used in summing these scores to produce aggregate scores. The top scoring sentences are used to create the summary itself.

#### 3.2.1 Trivial Systems

Two of our content selection systems are trivial systems. One system scores sentences based on their length relative to the longest sentence in the document cluster, favoring longer sentences. The second scores sentences based on their position in their document, giving the first sentence in each document a score of 1.0 and the last sentence in each document a score of 0.0.

#### 3.2.2 TF-IDF Scoring

For our first non-trivial content selection system, we calculated the average tf-idf score of every non-stop word in the sentence. TF-IDF scores were calculated using the Reuters-21578, Distribution 1.0 Corpus of news articles, as incorporated in the NLTK, as our background corpus. Sentence scores are scaled such that the sentence with the highest average TF-IDF was given a score of 1.0.

#### 3.2.3 Simple Graph Based Scoring

Next, we implemented a simple graph based metric, in which a dense undirected graph of sentences is constructed with edge weights set to the cosine similarity of the sentences. The most connected sentence is iteratively selected, and the weights of the edges of the previously selected sentences are set to be negative to discourage redundancy. The first pulled sentence is given a score of 1.0, with scores incrementally decreasing to 0.0 for the sentence pulled last.

#### 3.2.4 Noun Phrase Clustering

Finally, we also implemented a noun phrase clustering algorithm where we used a custom coreference resolution system to resolve pronouns to their most likely antecedent, and then compared each sentence's noun phrases to each other sentence. NP-clustering selection was ordered by the number of matches each sentence had to every other sentence.

### 3.3 Information Ordering

In this deliverable we did not attempt any global information ordering, leaving the summary sentences in the order of importance as determined by each content selection system (most similarities for graph similarity, highest clusters for noun phrase clustering, etc).

### 3.4 Content Realization

Our realization was simple, the highest ranking sentences were added into the summary with only extra white space and newlines removed. We iterated through the aggregate scored sentences in order, adding each sentence if there was room left within the one hundred word length limit.

### 3.5 Best Technique

The best system turned out to be the noun phrase clustering one, followed closely by an "equal weight" between all the systems.

## 4 Results

Our best results are shown below in Table 1. It turned out to be a close race, with our NP-Clustering-only solution (Table 1) winning over an equal weight system (Table 2) by a slim margin. We also ran the tf-idf by itself (Table 4), as well as Simple Graph only (Table 3). We found that our techniques generally worked the best when working together, but the NP Clustering by itself seemed to be the exception.

Table 1: NP Clustering

Rouge Technique	Recall	Precision	F-Score
ROUGE-1	0.23391	0.28553	0.25522
ROUGE-2	0.05736	0.07053	0.06272
ROUGE-3	0.01612	0.01969	0.01758
ROUGE-4	0.00533	0.00657	0.00584

Table 2: All Techniques

Rouge Technique	Recall	Precision	F-Score
ROUGE1	0.23336	0.28628	0.25516
ROUGE2	0.05708	0.07044	0.06251
ROUGE3	0.01612	0.01969	0.01758
ROUGE4	0.00533	0.00657	0.00584

Table 3: Simple Graph

Rouge Technique	Recall	Precision	F-Score
ROUGE-1	0.19379	0.25550	0.21845
ROUGE-2	0.04473	0.05859	0.05033
ROUGE-3	0.01170	0.01505	0.01305
ROUGE-4	0.00362	0.00453	0.00400

Table 4: tf-idf

Rouge Technique	Recall	Precision	F-Score
ROUGE-1	0.15341	0.20846	0.17522
ROUGE-2	0.03014	0.04037	0.03426
ROUGE-3	0.00746	0.01038	0.00863
ROUGE-4	0.00242	0.00329	0.00278

### 4.1 Error Analysis

There are many things to still tweak with both noun phrase clustering and graph similarity. For example, we could do better normalization with other co-referents rather than just pronouns. We will be experimenting with these techniques further.

## 5 Discussion

TODO: fill out

## 6 Conclusion

Conclusions have been made as can be seen from the following Nenkova, Radev, and Jones (Nenkova et al., 2007) (Jones, 2007) (Radev et al., 2001). We used co-referenced based off ideas from (Cardie and Wagstaff, 1999).

## References

- C. Cardie and K. Wagstaff. 1999. Noun phrase coreference as clustering. *Proceedings of the Joint Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*, 1001:82–89.
- Karen Spärck Jones. 2007. Automatic summarising: The state of the art. *Information Processing Management*, 43(6):1449 – 1481. Text Summarization.
- Ani Nenkova, Rebecca Passonneau, and Kathleen McKeown. 2007. The pyramid method: Incorporating human content selection variation in summarization evaluation. *ACM Transactions on Speech and Language Processing (TSLP)*, 4(2):4.
- Dragomir R Radev, Sasha Blair-Goldensohn, and Zhu Zhang. 2001. Experiments in single and multi-document summarization using mead. *Ann Arbor*, 1001:48109.