

# The University of Michigan at DUC 2004

Güneş Erkan<sup>1</sup>, Dragomir R. Radev<sup>2,1</sup>

<sup>1</sup>Department of EECS, <sup>2</sup>School of Information  
University of Michigan  
{gerkan,radev}@umich.edu

## Abstract

We present the results of Michigan's participation in DUC 2004. Our system, MEAD, ranked as one of the top systems in four of the five tasks. We introduce our new feature, LexPageRank, a new measure of sentence centrality inspired by the prestige concept in social networks. LexPageRank gave promising results in multi-document summarization. Our approach for Task 5, biographical summarization, was simplistic, yet successful. We used regular expression matching to boost up the scores of the sentences that are likely to contain biographical information patterns.

## 1 Introduction

The year 2004 marked the fourth time the University of Michigan's CLAIR (Computational Linguistics And Information Retrieval) group participated in the DUC evaluation. We entered our system, MEAD (Radev et al., 2001), in all of the five tasks, submitting 3 runs for each of the first 4 tasks and only one run for Task 5. The performance of our runs (peers numbered 135-147 in the official results) was quite good - our system ranked one of the top with respect to most of the metrics used in 4 of the 5 tasks.

In this report we will describe our general approach to the different tasks, paying particular attention to the ways in which we adapted our existing extractive summarizer, MEAD, to perform different tasks.

## 2 The DUC 2004 evaluation

The DUC 2004 evaluation consisted of 5 tasks. Tasks 1 and 2 were essentially the same as in last year's evaluation. Tasks 3 and 4 were also similar to Tasks 1 and 2, respectively, with the exception that the documents used in summarization are either automatic or manual English translations of Arabic documents.

The goal of Task 1 was to produce a 75-byte single-document summaries (headlines) for each of the documents in the given 50 TDT clusters. In Task 2, same set of clusters were used to produce a 665-byte multi-document short summary for each cluster.

The goals of Tasks 3 and 4 were essentially the same as Tasks 1 and 2, respectively, except that they focused on cross-lingual summaries of 25 Arabic TDT clusters. There were 3 subtasks for each of the Tasks 3 and 4 according to the data sets being used. In subtasks 3a and 4a, automatic English translations of the Arabic clusters were used. In subtasks 3b and 4b, manual translations were provided instead. Finally, in subtasks 3c and 4c, automatic translations plus some original English documents that were relevant to each cluster were given as the input.

Task 5 involved a newly introduced summarization concept for DUC, which is known as "biographical summarization". Given a news cluster of approximately 10 documents, and a question of the form "Who is X", where X is the name of a person, the task was to produce a 665-byte multi-document summary of the cluster that would respond to the question. 50 TREC news clusters, each of which focused on events about a different person, were used in this task.

## 3 Evaluation measures

In DUC 2004, an automatic evaluation metric for summarization, ROUGE<sup>1</sup>, was used for the first time. ROUGE is a recall-based metric for fixed-length summaries which is based on n-gram co-occurrence. It reports separate scores for 1, 2, 3, and 4-gram, and also for longest common sub-sequence co-occurrences. Tasks 1-4 were evaluated solely by means of ROUGE.

Task 5 summaries were evaluated manually for quality, coverage, and responsiveness to the question. Here is the complete list of metrics used in the evaluation of Task 5:

---

<sup>1</sup><http://www.isi.edu/~cyl/ROUGE>

### Quality metrics:

**Answers to 7 quality questions:** The assessors answered 7 questions (see Figure 1 ) for each summary. These were multiple choice questions with answers numbered from 1 to 5, where 1 was the best answer, 5 was the worst answer.

### Coverage metrics:

**Number of peer units:** Number of rough sentences in a summary.

**Number of marked peer units:** Number of peer units that the assessor felt expressed at least some of the meaning of the model.

**Number of unmarked peer units:** Number of peer units that the assessor felt did not express any of the meaning of the model.

**Fraction of unmarked peer units at least related to the model's subject:** The fraction of the number of peer units which did not overlap at all in meaning with any model unit, but was at least related to the subject of the model.

**Number of model units:** The number of roughly elementary discourse units (e.g., clauses etc) in the model

**Mean coverage:** The assessor judges the coverage by the peer summary of each unit in the model. This is the mean of those coverage scores.

**Median coverage:** Median of the per-model-unit coverage scores.

**Sample std of coverage scores:** Sample standard deviation of the per-model-unit coverage scores.

### Responsiveness metrics:

**Responsiveness score:** An integer grade between [0,4], where 0 is the worst and 4 is the best score, indicating how responsive the summary is to the question relative to the other summaries.

Q1: Does the summary build from sentence to sentence to a coherent body of information about the topic? Q2: If you were editing the summary to make it more concise and to the point, how much useless, confusing or repetitive text would you remove from the existing summary? Q3: To what degree does the summary say the same thing over again? Q4: How much trouble did you have identifying the referents of noun phrases in this summary? Are there nouns, pronouns or personal names that are not well-specified? For example, a person is mentioned and it is not clear what his role in the story is, or any other entity that is referenced but its identity and relation with the story remains unclear. Q5: To what degree do you think the entities (person/thing/event/place/...) were re-mentioned in an overly explicit way, so that readability was impaired? For example, a pronoun could have been used instead of a lengthy description, or a shorter description would have been more appropriate? Q6: Are there any obviously ungrammatical sentences, e.g., missing components, unrelated fragments or any other grammar-related problem that makes the text difficult to read? Q7: Are there any datelines, system-internal formatting or capitalization errors that can make the reading of the summary difficult?
--

Figure 1: The seven *quality* questions used in DUC 2004.

Centroid 1 Position 1 LengthCutoff 9 SimWithFirst 2 LexPageRank 1 mmr-reranker-word.pl 0.5 MEAD-cosine enidf
--

Figure 2: Sample MEAD policy.

## 4 Our system

We used the latest version of the MEAD system<sup>2</sup> augmented with a number of new rerankers. For a detailed discussion of MEAD, we refer the reader to (Radev et al., 2001). Suffice it to say that MEAD is an extractive summarization environment based on a three-step architecture. During the first step, *the feature extractor*, each sentence in the input document (or cluster of documents) is converted into a feature vector using features such as Position, Centroid, Length, OverlapWithFirst, etc. Second, the feature vector is converted to a scalar value using the *combiner*. At the last stage known as the *reranker*, the scores for sentences included in related pairs are adjusted upwards or downwards based on the type of relation between the sentences in the pair. Generally speaking, a *negative* relation exists between sentences that overlap in content (e.g., sentence pairs exhibiting subsumption or paraphrase) and therefore the presence of one of them in the summary should suppress the other one, while sentence pairs are related *positively* if the presence of one of them requires the presence of the other (e.g., due to an anaphoric relationship between them). The third stage of the MEAD architecture is based on custom *rerankers* which adjust the sentence scores assigned by the first and second stages. We used several rerankers in our experiments. Some of them (e.g., Maximal Marginal Relevance, MMR), are based on work by others (Carbonell and Goldstein, 1998) while others are based on our CST theory (Radev, 2000).

A MEAD policy is a combination of three components: (a) the command lines for all features, (b) the formula for converting the feature vector to a scalar, and (c) the command line for the reranker. A sample policy might be the one shown in Figure 2. This example indicates the four features used (Centroid, Position, LengthCutoff, and SimWithFirst), their relative weights (except for LengthCutoff where the number 9 indicates the threshold for selecting a sentence based on length), and the reranker (in this example, word-based MMR with a similarity threshold computed as the cosine between two sentences).

### 4.1 Features

We used four of the old features of MEAD in all of the tasks. These are Centroid, Position, LengthCutoff, and SimWithFirst. Additionally, we introduced two new features for this year's DUC evaluation.

<sup>2</sup><http://www.summarization.com>

SNo	ID	Text
1	d1s1	Iraqi Vice President Taha Yassin Ramadan announced today, Sunday, that Iraq refuses to back down from its decision to stop cooperating with disarmament inspectors before its demands are met.
2	d2s1	Iraqi Vice president Taha Yassin Ramadan announced today, Thursday, that Iraq rejects cooperating with the United Nations except on the issue of lifting the blockade imposed upon it since the year 1990.
3	d2s2	Ramadan told reporters in Baghdad that "Iraq cannot deal positively with whoever represents the Security Council unless there was a clear stance on the issue of lifting the blockade off of it.
4	d2s3	Baghdad had decided late last October to completely cease cooperating with the inspectors of the United Nations Special Commission (UNSCOM), in charge of disarming Iraq's weapons, and whose work became very limited since the fifth of August, and announced it will not resume its cooperation with the Commission even if it were subjected to a military operation.
5	d3s1	The Russian Foreign Minister, Igor Ivanov, warned today, Wednesday, against using force against Iraq, which will destroy, according to him, seven years of difficult diplomatic work and will complicate the regional situation in the area.
6	d3s2	Ivanov contended that carrying out air strikes against Iraq, who refuses to cooperate with the United Nations inspectors, "will end the tremendous work achieved by the international group during the past seven years and will complicate the situation in the region."
7	d3s3	Nevertheless, Ivanov stressed that Baghdad must resume working with the Special Commission in charge of disarming the Iraqi weapons of mass destruction (UNSCOM).
8	d4s1	The Special Representative of the United Nations Secretary-General in Baghdad, Prakash Shah, announced today, Wednesday, after meeting with the Iraqi Deputy Prime Minister, Tariq Aziz, that Iraq refuses to back down from its decision to cut off cooperation with the disarmament inspectors.
9	d5s1	British Prime Minister Tony Blair said today, Sunday, that the crisis between the international community and Iraq "did not end" and that Britain is still "ready, prepared, and able to strike Iraq."
10	d5s2	In a gathering with the press held at the Prime Minister's office, Blair contended that the crisis with Iraq "will not end until Iraq has absolutely and unconditionally respected its commitments" towards the United Nations.
11	d5s3	A spokesman for Tony Blair had indicated that the British Prime Minister gave permission to British Air Force Tornado planes stationed in Kuwait to join the aerial bombardment against Iraq.

	1	2	3	4	5	6	7	8	9	10	11
1	1.00	0.45	0.02	0.17	0.03	0.22	0.03	0.28	0.06	0.06	0.00
2	0.45	1.00	0.16	0.27	0.03	0.19	0.03	0.21	0.03	0.15	0.00
3	0.02	0.16	1.00	0.03	0.00	0.01	0.03	0.04	0.00	0.01	0.00
4	0.17	0.27	0.03	1.00	0.01	0.16	0.28	0.17	0.00	0.09	0.01
5	0.03	0.03	0.00	0.01	1.00	0.29	0.05	0.15	0.20	0.04	0.18
6	0.22	0.19	0.01	0.16	0.29	1.00	0.05	0.29	0.04	0.20	0.03
7	0.03	0.03	0.03	0.28	0.05	0.05	1.00	0.06	0.00	0.00	0.01
8	0.28	0.21	0.04	0.17	0.15	0.29	0.06	1.00	0.25	0.20	0.17
9	0.06	0.03	0.00	0.00	0.20	0.04	0.00	0.25	1.00	0.26	0.38
10	0.06	0.15	0.01	0.09	0.04	0.20	0.00	0.20	0.26	1.00	0.12
11	0.00	0.00	0.00	0.01	0.18	0.03	0.01	0.17	0.38	0.12	1.00

Figure 3: Intra-sentence cosine similarities in a subset of cluster d1003t from DUC 2004.

#### 4.1.1 LexPageRank feature

Extractive summarization can be viewed as choosing the most *central* sentences in a cluster that give the necessary and enough amount of information related to the main theme of the cluster. In the case of centroid-based summarization (Radev et al., 2000), *centrality* is defined in terms of the centroid of the cluster, i.e. the sentences that contain the highest ranked words in the centroid of the cluster are considered as central to the topic.

We propose a new measure of sentence centrality inspired by the concept of *prestige* in social networks and its application in the Web. PageRank (Page et al., 1998) is a method proposed for assigning a prestige score to each page in the web independent of a specific query. In PageRank, the score of a page is determined depending on the number of pages that link to that page as well as the individual scores of the linking pages. This is achieved

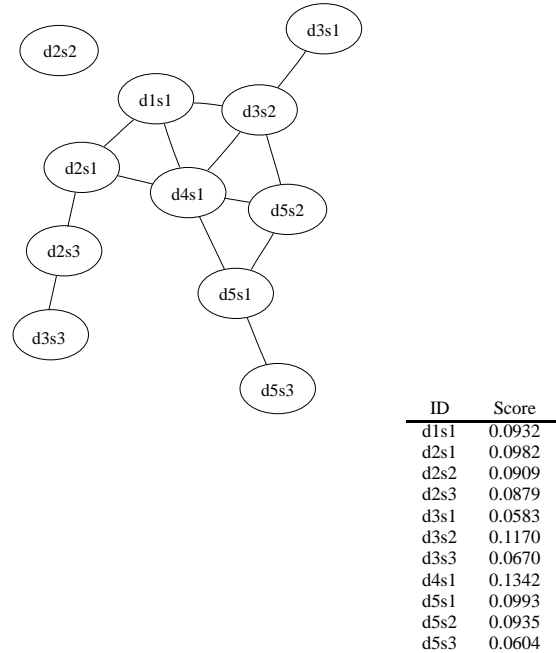


Figure 4: Similarity graph and corresponding LexPageRank scores for the cluster in Figure 3. Sentence d4s1 is the most central page (score = 0.1342).

by forming the normalized adjacency matrix of the network of pages and using a power method to converge to the principal eigenvector of this matrix. In the context of summarization, we recursively define the prestige of a sentence in terms of the number of sentences that it is similar to, and the corresponding prestiges of these similar sentences. We use cosine as the similarity measure and define the cosine adjacency matrix of a cluster as the binary connectivity matrix,  $M$ , where  $M(u, v) = 1$  if the cosine similarity between the sentence  $u$  and the sentence  $v$  is above a certain threshold. Figure 3 shows a subset of a cluster used in Task 4b of DUC 2004, and the cosine similarity between each pair of sentences. Sentence ID  $dXsY$  indicates the  $Y$ th sentence in the  $X$ th document. Figure 4 shows the graph that corresponds to the similarity matrix with threshold 0.2. The LexPageRank scores in the same figure are computed by an approximation of the principal eigenvector of the adjacency matrix for the graph. Unlike the original PageRank method, the graph is undirected since cosine similarity is symmetric.

We used LexPageRank feature in our priority 3 submission for Task 2, and priority 1 submission for Task 4. The cosine threshold used was 0.3. Although we did not experiment with enough number of different thresholds and feature weights, the results are better according to many of the metrics compared to our submissions without LexPageRank (Table 5). A more detailed performance analysis for LexPageRank will be presented in an

extended version of this paper.

#### 4.1.2 QueryPhraseMatch feature

Task 5 needed a very different treatment from other tasks since it involved producing biographical summaries, which are quite different from general purpose summaries. In order to do this, we searched for a way of giving higher scores to the sentences that gives information about the person in the question. There could be two ways of incorporating such a ranking of sentences into MEAD; either introducing a new feature that gives higher scores to such sentences, or a new reranker that modifies the rankings of the sentences in a general purpose summary. We chose the first one, and introduced a new feature called QueryPhraseMatch.

In order to compute the QueryPhraseMatch feature for a sentence, we look for certain regular expressions in the sentence and increase the feature value for each expression that matches the sentence. Note that QueryPhraseMatch is a general purpose feature, not designed solely for the purpose of biographical summaries, and can be used to look for any kind of regular expressions in a text. In the case of a biographical summary, these should be the expressions that often occur in sentences that describe a person (e.g. an occurrence of person’s name, relative clauses, etc.). We looked at some biographical texts and collected regular expressions that appear in these texts in common. A subset of these expressions are shown in Figure 5. Every expression has a weight that determines by how much we will increase the score of a sentence that matches the expression.

Due to lack of time, the regular expressions and corresponding weights were determined totally empirically and depending on subjective judgements. However, the results are promising as we discuss in Section 5.

Expression	Weight
X	0.25
X grew up	1
X attended	1
X (turns turned) [1-9][0-9]?	1
X, (an?) the who whom whose) [\w ]*[.,]	1.5
X, [1-9][0-9]?(\. years)	1.5
X began	0.35
X (lives lived)	0.5
X made	0.5

Figure 5: Some regular expressions that are searched to determine the sentences that describe a person. ‘X’ is unified with the person name at run time.

## 4.2 Training

We participated in the headline extraction tasks (Tasks 1 and 3) for the first time in DUC 2004. We used MEAD to extract only one sentence from a document as the headline of that document. No real training was done for these tasks. However, we observed that setting the weight of

the Position feature high results in higher ROUGE scores for DUC 2003 data, meaning that the sentences that occur at the beginning of a document are more likely to get higher ROUGE scores with respect to the manual headlines. This issue is discussed in Section 5. Table 1 shows the feature weights we used in our submissions for Tasks 1 and 3.

Feature Weights			Task 1	Task 3
Centroid	Position	Length	Priority	Priority
1.0	1.0	7	1	3
1.0	1.5	7	2	1
1.0	10.0	7	3	2

Table 1: Feature weights used in Tasks 1 and 3.

To train MEAD for Tasks 2 and 4, we used the DUC 2003 data set. We split the data into two parts as training and devtest. We fixed our reranker to be the MMR-reranker since initial experiments showed that it outperforms the other rerankers in almost any policy as far as ROUGE scores are concerned. All of the experiment results shown in this paper were derived from MEAD policies that use MMR-reranker.

We made a local search with random restarts on the feature weights space. For each set of feature weights, we ran MEAD on the test data, and tried to maximize the ROUGE-1 (unigram) and ROUGE-W (weighted longest common subsequence) scores. The search step size was decreased as we got better ROUGE scores. Table 2 shows the ROUGE scores we got on the devtest data with the best policies observed on the training data. The highlighted rows are the two policies we used in our submissions. One of them gave the highest ROUGE-W score while the other was the one which gave the highest ROUGE-1 score among the policies which were substantially different than the first one.

Centroid	Position	Length	SimWithFirst	ROUGE-1	ROUGE-W
1.5	1.0	9	2.0	0.32256	0.11712
1.5	1.0	9	3.0	0.32256	0.11712
2.0	1.0	9	2.0	0.32750	0.11875
3.0	1.0	9	2.0	0.33311	0.11847
3.0	1.0	15	2.0	0.34171	0.12132
3.0	1.0	20	2.0	0.33566	0.11906
<b>2.0</b>	<b>1.0</b>	<b>9</b>	<b>4.0</b>	<b>0.34361</b>	<b>0.12511</b>
3.0	1.0	9	4.0	0.33926	0.12301
4.0	1.0	9	4.0	0.34560	0.12171
3.0	1.0	9	5.0	0.34623	0.12174
2.0	1.0	15	4.0	0.34470	0.12223
2.0	1.0	9	5.0	0.34707	0.12304
<b>3.0</b>	<b>1.0</b>	<b>15</b>	<b>3.0</b>	<b>0.34684</b>	<b>0.12197</b>
3.0	1.0	15	4.0	0.34606	0.12170

Table 2: Best policies and corresponding ROUGE scores for DUC 2003 devtest data

Due to lack of time, we were not able to experiment

with the weight of our new LexPageRank feature. Instead, we picked the policies that performed the best without LexPageRank and reran them with LexPageRank using a constant weight. Table 3 shows the ROUGE scores for these policies. The highlighted policy is the one we used in our submissions for Tasks 2 and 4.

Centroid	Position	Length	SimWithFirst	LexPageRank	ROUGE-1	ROUGE-W
2.0	1.0	9	4.0	2.0	0.34129	0.12245
3.0	1.0	9	4.0	2.0	0.34051	0.12182
1.0	1.0	10	4.0	2.0	0.33662	0.12059
1.0	1.0	11	5.0	2.0	0.33806	0.12137
4.0	1.0	9	4.0	2.0	0.34513	0.12233
3.0	1.0	9	5.0	2.0	0.34450	0.12322
2.0	1.0	15	4.0	2.0	0.32436	0.11532
<b>2.0</b>	<b>1.0</b>	<b>9</b>	<b>5.0</b>	<b>2.0</b>	<b>0.34450</b>	<b>0.12370</b>
3.0	1.0	15	3.0	2.0	0.33922	0.12090
3.0	1.0	15	4.0	2.0	0.34467	0.12280

Table 3: Policies with LexPageRank feature and corresponding ROUGE scores for DUC 2003 devtest data

No training was done for Task 5. MEAD was used with default feature weights. However, the weights of the regular expressions incorporated into the QueryPhraseMatch feature were high enough to override default MEAD features.

## 5 Results

Table 5 shows the rankings of our submissions for Tasks 1–4 with respect to several ROUGE metrics among all of the submissions. We chose to separate the results for the subtasks of Tasks 3 and 4 since the data sets for these subtasks were different from each other.

Peer Code	Task	Priority	Rankings						Total Submissions
			ROUGE-1	ROUGE-2	ROUGE-3	ROUGE-4	ROUGE-L	ROUGE-W	
135	1	1	7	4	4	4	5	5	39
136	1	2	5	3	3	3	4	3	39
137	1	3	2	2	1	1	2	2	39
138	2	1	18	16	16	13	27	27	34
139	2	2	21	18	17	17	29	29	34
140	2	3	15	20	19	16	26	26	34
141	3a	NA	5	2	1	1	2	2	11
142	3b	NA	5	1	1	1	4	3	10
144	4a	NA	1	2	1	1	6	6	11
145	4b	NA	3	1	1	1	7	7	11
146	4c	NA	1	2	2	2	4	4	6

Table 5: Official rankings of our submissions for Tasks 1–4 with respect to ROUGE metrics.

On Task 1, our best system (peer code: 137) was the priority-3 submission for this task with the highest weight for Position feature among our tasks (Table 1), and supports our observations in the experiments that sentences

that occur near the beginning of a document are more likely to get higher ROUGE scores. This also explains the surprisingly successful performance of the baseline system on Task 1, which simply extracts the first sentence of a document.

Task 2 was our most unsuccessful task, where our system showed an average performance among all of the submissions. Our submission that included the LexPageRank feature (peer code: 140) achieved the best ROUGE-1 score among our three submissions for this task, which is a promising result for this new feature.

On Task 3, almost all systems participated in the evaluation got the best score on Task 3b since the data set was composed of manual translations, which are expected to be more similar to model human summaries compared to automatic translations. Our best score came from the best policy on Task 1, which we again used here for Task 3b. Since only two systems participated in Task 3c, we did not include the ranking of our submission for this task in Table 5. Furthermore, our submission can essentially be compared to the submissions for Task 3a since we did not use the relevant documents provided by NIST for Task 3c. The difference between the scores of our submissions for Task 3a and Task 3c is due to the difference between the weights of the Position feature in two policies. As in Task 1, we got higher ROUGE scores with the higher Position weight.

Our results on Task 4 are also very successful. The results we got using LexPageRank feature (peer code: 145) are as successful as our other submissions for this task, which is again a promising fact.

There are several metrics for evaluating Task 5. However, these metrics are not defined as an overall evaluation of system performance but rather as an assessment of each individual summary. Table 4 shows our own analysis of average performance of each system. Numbers in the parentheses show the number of clusters in which the system was one of the best performers among all systems (excluding manual summaries) considering the metric of that column only. First seven columns show the average scores of the systems for each quality question in Figure 1. A lower number shows a better system. Since our system produces extractive summaries, its performance is one of the best in questions about grammaticality and formatting (e.g. Q6 and Q7). However, it performs worse in questions about the semantic structure of the summary (e.g. Q3, Q4, and Q5). Next column shows the mean of the mean coverage scores that a system got for each cluster. Our system ranked third overall, and ranked best in 4 of the 50 clusters. Another metric for Task 5 is the fraction of the unmarked peer units, i.e. the units which do not overlap at all in meaning with any model unit, that are somewhat related to the model’s subject. We took the average of these fractions of a peer system for all clus-

Peer code	Q1	Q2	Q3	Q4	Q5	Q6	Q7	Mean coverage	Unmarked but related units (%)	Responsiveness
A	1.00	1.06	1.00	1.00	1.00	1.11	1.06	0.477	100	3.60
B	1.06	1.06	1.00	1.06	1.00	1.06	1.00	0.554	100	3.92
C	1.24	1.36	1.28	1.04	1.16	1.40	1.12	0.350	86	2.84
D	1.11	1.22	1.00	1.06	1.17	1.33	1.00	0.469	96	3.60
E	1.39	1.33	1.00	1.11	1.00	1.33	1.22	0.468	90	3.32
F	1.24	1.24	1.06	1.00	1.06	1.12	1.18	0.517	84	3.64
G	1.17	1.00	1.00	1.00	1.00	1.28	1.06	0.537	94	3.48
H	1.17	1.22	1.11	1.00	1.06	1.06	1.00	0.438	96	3.52
5	1.62 (38)	2.20 (28)	1.44 (37)	1.44 (37)	1.10 (47)	1.40 (33)	1.82 (16)	0.190 (2)	62	1.52 (16)
16	3.12 (8)	2.70 (12)	1.58 (34)	2.04 (20)	1.70 (28)	1.36 (37)	1.46 (36)	0.158 (1)	69	1.30 (7)
24	3.86 (1)	3.40 (2)	1.44 (36)	2.98 (7)	1.30 (38)	1.88 (25)	1.88 (20)	0.213 (8)	55	1.16 (6)
30	2.90 (6)	2.42 (20)	1.38 (39)	2.00 (24)	1.46 (33)	1.30 (39)	1.30 (40)	0.199 (3)	72	1.42 (13)
43	3.78 (1)	2.94 (13)	1.30 (38)	2.84 (8)	1.46 (32)	1.76 (28)	1.52 (40)	0.198 (6)	70	1.30 (14)
49	3.34 (5)	2.88 (9)	1.54 (36)	2.12 (19)	1.62 (33)	1.78 (25)	1.36 (38)	0.206 (3)	78	1.40 (12)
62	3.62 (3)	3.12 (10)	1.98 (24)	1.90 (24)	2.14 (18)	1.84 (21)	1.56 (36)	0.200 (3)	49	1.64 (17)
71	3.32 (2)	2.96 (3)	1.84 (25)	1.44 (36)	2.42 (12)	1.64 (26)	1.66 (29)	0.214 (2)	64	1.50 (13)
86	3.86 (3)	3.22 (5)	1.16 (44)	2.96 (10)	1.18 (41)	2.22 (16)	1.68 (32)	0.145 (2)	51	1.06 (9)
96	3.30 (10)	3.06 (8)	1.60 (29)	1.82 (28)	1.66 (30)	1.68 (26)	3.54 (3)	0.216 (6)	60	1.44 (11)
109	3.14 (6)	2.76 (11)	1.86 (26)	1.74 (28)	1.96 (25)	1.42 (35)	1.54 (38)	0.241 (7)	76	1.76 (19)
116	4.52 (0)	4.04 (1)	1.48 (35)	3.40 (3)	1.76 (24)	2.14 (17)	2.14 (15)	0.173 (7)	65	1.00 (8)
122	2.94 (11)	2.32 (27)	1.68 (30)	1.24 (43)	2.08 (18)	1.22 (43)	1.36 (42)	0.184 (1)	74	1.26 (7)
125	3.52 (4)	3.22 (6)	1.62 (28)	2.12 (23)	1.94 (21)	2.76 (8)	2.18 (23)	0.189 (4)	71	1.40 (13)
<b>147</b>	<b>3.18 (8)</b>	<b>2.64 (14)</b>	<b>1.84 (24)</b>	<b>1.92 (24)</b>	<b>1.80 (25)</b>	<b>1.24 (39)</b>	<b>1.40 (39)</b>	<b>0.215 (4)</b>	<b>69</b>	<b>1.54 (19)</b>

Table 4: An overall analysis of official results for Task 5.

ters, weighting by the number of unmarked units of the system for each cluster. As seen in Table 4, about 70% of the unmarked units of our summaries were at least related to the subject. The last column shows the average responsiveness scores. Our system ranked third overall while ranking best in 19 of 50 clusters.

## 6 Post-DUC experiments

After the official evaluations, we carefully reimplemented LexPageRank to see how this new method performs compared to centroid-based summarization as well as to other DUC participant systems. We ran MEAD with several policies with different feature weights and combinations of features. However, we did not use Centroid and LexPageRank features in a same policy to get a better comparison of two methods. We fixed Length cutoff at 9, and the weight of the Position feature at 1 in all of the policies. We did not try a weight higher than 2.0 for any of the features since our earlier observations on MEAD showed that too high feature weights results in poor summaries.

Table 6 and Table 7 show the ROUGE scores we have got in the experiments with using LexPageRank and Centroid in Tasks 2 and 4, respectively, sorted by ROUGE-1 scores. ‘lprXTY’ indicates a policy in which the weight for LexPageRank is  $X$  and  $Y$  is used as threshold. ‘CX’

shows a policy with Centroid weight  $X$ . We also include two baselines for each data set. ‘random’ indicates a method where we have picked random sentences from the cluster to produce a summary. We have performed five random runs for each data set. The results in the tables are for the median runs. Second baseline, shown as ‘lead-based’ in the tables, is using only the Position feature without any centrality method. This is tantamount to producing lead-based summaries, which is a widely used and very challenging baseline in the text summarization community (Brandow et al., 1995).

The top scores we have got in all data sets come from our new method, LexPageRank. The results provide strong evidence that LexPageRank is better than Centroid in multi-document generic text summarization. Another interesting observation in the results is the effect of threshold. Most of the top ROUGE scores belong to the runs with the threshold 0.1, and the runs with threshold 0.3 are worse than the others most of the time. This is due to the information loss in the similarity graphs as we move to higher thresholds since a higher threshold gives us a sparser similarity graph. The results suggest that 0.1 is a suitable threshold for LexPageRank compared to higher numbers like 0.3 which we used in the official runs.

As a comparison with the other summarization sys-

tems, we present the official scores for the top five DUC 2004 participants and the human summaries in Table 8 and Table 9 for Tasks 2 and 4, respectively. Our top few results for each task are either better than or statistically indifferent from the best system in the official runs considering the 95% confidence interval.

Policy Code	ROUGE-1 (unigram)	ROUGE-2 (bigram)	ROUGE-W (LCS)
lpr2T0.1	0.38079	0.08971	0.12984
lpr1.5T0.1	0.37873	0.09068	0.13032
lpr0.5T0.1	0.37842	0.08972	0.13121
lpr1T0.1	0.37700	0.09174	0.13096
C0.5	0.37672	0.09233	0.13230
lpr1T0.2	0.37667	0.09115	0.13234
lpr0.5T0.2	0.37482	0.09160	0.13220
C1	0.37464	0.09210	0.13071
lpr1T0.3	0.37448	0.08767	0.13302
lpr0.5T0.3	0.37362	0.08981	0.13173
lpr1.5T0.2	0.37058	0.08658	0.12965
C1.5	0.36885	0.08765	0.12747
lead-based	0.36859	0.08669	0.13196
lpr1.5T0.3	0.36849	0.08455	0.13111
lpr2T0.3	0.36737	0.08182	0.13040
lpr2T0.2	0.36737	0.08264	0.12891
C2	0.36710	0.08696	0.12682
random	0.32381	0.05285	0.11623

Table 6: Results for Task 2

## 7 Acknowledgments

This work was partially supported by the National Science Foundation under grant 0329043 “Probabilistic and link-based Methods for Exploiting Very Large Textual Repositories” administered through the IDM program. All opinions, findings, conclusions, and recommendations in this paper are made by the authors and do not necessarily reflect the views of the National Science Foundation.

We would like to thank the members of the CLAIR group at Michigan and in particular Jahna Otterbacher and Siwei Shen for their assistance with this project.

## References

- Ron Brandow, Karl Mitze, and Lisa F. Rau. 1995. Automatic condensation of electronic publications by sentence selection. *Information Processing and Management*, 31(5):675–685.
- Jaime G. Carbonell and Jade Goldstein. 1998. The use of MMR, diversity-based reranking for reordering documents and producing summaries. In *Research and Development in Information Retrieval*, pages 335–336.
- L. Page, S. Brin, R. Motwani, and T. Winograd. 1998. The pagerank citation ranking: Bringing order to the web. *Technical report, Stanford University, Stanford, CA*.
- Dragomir R. Radev, Hongyan Jing, and Malgorzata Budzikowska. 2000. Centroid-based summarization

Policy Code	ROUGE-1 (unigram)	ROUGE-2 (bigram)	ROUGE-W (LCS)
Task 4a			
lpr1.5T0.1	0.39997	0.11030	0.12427
lpr1.5T0.2	0.39970	0.11508	0.12422
lpr2T0.2	0.39954	0.11417	0.12468
lpr2T0.1	0.39809	0.11033	0.12357
lpr1T0.2	0.39614	0.11266	0.12350
lpr0.5T0.1	0.39369	0.10665	0.12287
lpr1T0.1	0.39312	0.10730	0.12274
C0.5	0.39013	0.10459	0.12202
lpr0.5T0.2	0.38899	0.10891	0.12200
lpr1T0.3	0.38777	0.10586	0.12157
lpr0.5T0.3	0.38667	0.10255	0.12244
lpr1.5T0.3	0.38251	0.10610	0.12039
C1	0.38181	0.10023	0.11909
lpr2T0.3	0.38096	0.10497	0.12001
C1.5	0.38074	0.09922	0.11804
C2	0.38001	0.09901	0.11772
lead-based	0.37880	0.09942	0.12218
random	0.35929	0.08121	0.11466
Task 4b			
lpr1.5T0.1	0.40639	0.12419	0.13445
lpr2T0.1	0.40529	0.12530	0.13346
C1.5	0.40344	0.12824	0.13023
C2	0.39997	0.12367	0.12873
lpr2T0.3	0.39859	0.11744	0.12924
lpr1.5T0.3	0.39858	0.11737	0.13044
lpr1.5T0.2	0.39819	0.12228	0.12989
lpr2T0.2	0.39763	0.12114	0.12924
lpr1T0.1	0.39552	0.12045	0.13304
lpr1T0.2	0.39492	0.12056	0.13061
C1	0.39388	0.12301	0.12805
lpr1T0.3	0.39053	0.11500	0.13044
lpr0.5T0.1	0.38374	0.11331	0.12954
lpr0.5T0.2	0.38201	0.11201	0.12757
C0.5	0.37601	0.11123	0.12605
lpr0.5T0.3	0.37525	0.11115	0.12898
random	0.37339	0.09225	0.12205
lead-based	0.35872	0.10241	0.12496

Table 7: Results for Task 4

of multiple documents: sentence extraction, utility-based evaluation, and user studies. In *ANLP/NAACL Workshop on Summarization*, Seattle, WA, April.

Dragomir Radev, Sasha Blair-Goldensohn, and Zhu Zhang. 2001. Experiments in single and multi-document summarization using MEAD. In *First Document Understanding Conference*, New Orleans, LA, September.

Dragomir Radev. 2000. A common theory of information fusion from multiple text sources, step one: Cross-document structure. In *Proceedings, 1st ACL SIGDIAL Workshop on Discourse and Dialogue*, Hong Kong, October.

Peer Code	ROUGE-1 (unigram)	95% Confidence Interval	ROUGE-2 (bigram)	95% Confidence Interval	ROUGE-W (LCS)	95% Confidence Interval
H	0.4183	[0.4019,0.4346]	0.1050	[0.0902,0.1198]	0.1480	[0.1409,0.1551]
F	0.4125	[0.3916,0.4333]	0.0899	[0.0771,0.1028]	0.1462	[0.1388,0.1536]
E	0.4104	[0.3882,0.4326]	0.0984	[0.0838,0.1130]	0.1435	[0.1347,0.1523]
D	0.4060	[0.3870,0.4249]	0.1065	[0.0947,0.1184]	0.1449	[0.1395,0.1503]
B	0.4043	[0.3795,0.4291]	0.0950	[0.0785,0.1114]	0.1447	[0.1347,0.1548]
A	0.3933	[0.3722,0.4143]	0.0896	[0.0792,0.1000]	0.1387	[0.1319,0.1454]
C	0.3904	[0.3715,0.4093]	0.0969	[0.0849,0.1089]	0.1381	[0.1317,0.1444]
G	0.3890	[0.3679,0.4101]	0.0860	[0.0721,0.0998]	0.1390	[0.1315,0.1465]
65	0.3822	[0.3708,0.3937]	0.0922	[0.0827,0.1016]	0.1333	[0.1290,0.1375]
104	0.3744	[0.3635,0.3854]	0.0855	[0.0770,0.0939]	0.1284	[0.1244,0.1324]
35	0.3743	[0.3615,0.3871]	0.0837	[0.0737,0.0936]	0.1338	[0.1291,0.1384]
19	0.3739	[0.3602,0.3875]	0.0803	[0.0712,0.0893]	0.1315	[0.1261,0.1368]
124	0.3706	[0.3578,0.3835]	0.0829	[0.0748,0.0909]	0.1293	[0.1252,0.1334]
.	.	.	.	.	.	.
2	0.3242	[0.3104,0.3380]	0.0641	[0.0545,0.0737]	0.1186	[0.1130,0.1242]
.	.	.	.	.	.	.
.	.	.	.	.	.	.

Table 8: Summary of official ROUGE scores for DUC 2004 Task 2. Peer codes: baseline(2), manual[A-H], and system submissions

Peer Code	ROUGE-1 (unigram)	95% Confidence Interval	ROUGE-2 (bigram)	95% Confidence Interval	ROUGE-W (LCS)	95% Confidence Interval
Y	0.44450	[0.42298,0.46602]	0.12815	[0.10965,0.14665]	0.14348	[0.13456,0.15240]
Z	0.43263	[0.40875,0.45651]	0.11953	[0.10186,0.13720]	0.14019	[0.13056,0.14982]
X	0.42925	[0.40680,0.45170]	0.12213	[0.10180,0.14246]	0.14147	[0.13361,0.14933]
W	0.41188	[0.38696,0.43680]	0.10609	[0.08905,0.12313]	0.13542	[0.12620,0.14464]
Task 4a						
144	0.38827	[0.36261,0.41393]	0.10109	[0.08680,0.11538]	0.11140	[0.10471,0.11809]
22	0.38654	[0.36352,0.40956]	0.09063	[0.07794,0.10332]	0.11621	[0.10980,0.12262]
107	0.38615	[0.35548,0.41682]	0.09851	[0.08225,0.11477]	0.11951	[0.11004,0.12898]
68	0.38156	[0.36420,0.39892]	0.09808	[0.08686,0.10930]	0.11888	[0.11255,0.12521]
40	0.37960	[0.35809,0.40111]	0.09408	[0.08367,0.10449]	0.12240	[0.11659,0.12821]
.	.	.	.	.	.	.
.	.	.	.	.	.	.
Task 4b						
23	0.41577	[0.39333,0.43821]	0.12828	[0.10994,0.14662]	0.13823	[0.12995,0.14651]
84	0.41012	[0.38543,0.43481]	0.12510	[0.10506,0.14514]	0.13574	[0.12638,0.14510]
145	0.40602	[0.36783,0.44421]	0.12833	[0.10375,0.15291]	0.12221	[0.11128,0.13314]
108	0.40059	[0.37002,0.43116]	0.12087	[0.10212,0.13962]	0.13011	[0.12029,0.13993]
69	0.39844	[0.37440,0.42248]	0.11395	[0.09885,0.12905]	0.12861	[0.12000,0.13722]
.	.	.	.	.	.	.
.	.	.	.	.	.	.

Table 9: Summary of official ROUGE scores for DUC 2004 Task 4. Peer codes: manual[W-Z], and system submissions