

# The Computational Linguistics Summarization Pilot Task

**Kokil Jaidka<sup>†</sup>, Muthu Kumar Chandrasekaran<sup>2</sup>, Rahul Jha<sup>3</sup>, Christopher Jones<sup>4</sup>**

**Min-Yen Kan<sup>2,5</sup>, Ankur Khanna<sup>2</sup>, Diego Molla-Aliod<sup>4</sup>, Dragomir R. Radev<sup>3</sup>, Francesco Ronzano<sup>6</sup> and Horacio Saggion<sup>6</sup>**

## Abstract

The Computational Linguistics (CL) Summarization Pilot Task was a pilot shared task to use citations to create summaries of scholarly research publications in the domain of computational linguistics. We describe the background for the task, corpus construction, evaluation methods for the pilot and survey the participating systems and their preliminary results. The experience gleaned from the pilot will assist in the proper organization of future shared task where difficulties with annotations and scale can be addressed.

## 1 Introduction

This paper describes the evolution and design of the CL pilot task<sup>1</sup> for the summarisation of computational linguistics research papers sampled from the Association of Computational Linguistics' (ACL) anthology. This task was concurrently publicized

with TAC 2014, although it is not formally affiliated with the same, and shares its basic structure and guidelines with the more formal TAC 2014 Biomedical Summarization (BiomedSumm) track. A development corpus of training “topics” from CL research papers was released, each comprising a main, cited paper along with associated citing papers. Participants were invited to enter their systems in a task-based evaluation, similar to BiomedSumm.

This paper will describe the participating systems and survey their results from the task-based evaluation.

## 2 Background

Recent works (Mohammad et al., 2009); (Abu-Jbara and Radev, 2011) in scientific document summarisation have used citation sentences or citances from citing papers to create a multi document summary of the reference paper (RP).

As proposed by (VU, 2010); (Hoang and Kan, 2010) the summarisation can be decomposed into finding the relevant documents, in this case, the citing papers (CPs), then selecting sentences from those papers that cite and justify the citation and finally generate the summary. To help tackle each

<sup>\*</sup> Authors appear in alphabetical order, with the exception of the coordinator of the task, whom is given the first authorship.

<sup>1</sup>This research is support in part by the Singapore National Research Foundation under its International Research Centre @ Singapore Funding Initiative and administered by the IDM Programme Office.

of these subproblems, we created gold standard datasets where human annotators identify the citations in each of about 10 randomly sampled citing papers for the RP.

A pilot study conducted in the information science domain indicated that most citations clearly refer to one or more specific aspects of the cited paper (Jaidka et al., 2013). For computational linguistics, we identified that the discourse facets being cited were usually the aim of the paper, methods followed and the results or implications of the work. Accordingly, we used a different set of discourse facets than BiomedSumm which suit CL papers better. Please note that this is a development corpus and only a training set is available for use now. Although, we plan to release a test set of documents for next year's evaluation, we report k fold cross-validated performance over the 10 documents for the systems registered for participation.

### 3 Corpus Construction

The CL community uses the ACL Anthology Reference Corpus (Bird et al., 2008) to evaluate and report performance of systems. To support further research in scientific document summarisation among the CL community and beyond we plan to build a manually annotated corpus using research papers sampled from the ACL Anthology. As first steps towards this goal we created an annotated development corpus by randomly sampling 10 documents from the ACL anthology.

We now describe our construction in detail. From the current, live ACL Anthology, there are approx 26K (exactly 25961) individual papers that are in the Anthology (as of 18 September 2014; including ones staged for publication but not actually published yet). These only include files that we have PDFs hosted (e.g., LREC is not represented as we don't hold these PDFs in the Anthology, just their metadata). This number is approximate as there are some files that are not publications (frontmatter, author indices) that are included. Culling all files before and including 2006, we get 13.8K (13838) publications, which include conference and journal articles. We randomized this list to remove any ordering affects. Starting from the top of the list, we used

Google search <sup>2</sup> on (18 September 2013) to search for the publication - first by using its Anthology ID as a query (e.g., "H89-2014.pdf") and not productive, re-queried by the title of the paper as a string (e.g., "Some Experiments with a Naive Bayes WSD System"). We look for a Scholar search result that shows # of citations. This was an approximation. We kept any paper with over 10 citations. Some papers had some similar versions that presented different citation rates; however, all of these were dropped anyways due to low citation rate. We vetted the citations from Google Scholar <sup>3</sup> for the citation spread being over 3 years as per citing papers' year of publication (as in Google Scholar). We did not attempt to check for publication years that Google Scholar doesn't report for some publications. We check only the earliest range manually to ensure that the citation is the correct one, as there are usually many examples of later citations. We also vetted that there were at least 10 of these citing sources available for download. Some candidates were dropped due to the few amount of available files that were freely downloadable from the Web. We ran a title search to find the paper in ACL Anthology Network <sup>4</sup> (AAN, February 2013 version). We inspected and listed the citing papers (incoming citations) Anthology ID, title and year where the citing papers were given in reverse chronological order. Note the citation count from Google / Google Scholar and AAN (Feb 2013 release) will differ substantially.

To report the final list of citing papers, we strived to provide at least 3 citing papers for each paper. To do so, we defined the following criteria in order or priority):

1. Non-list citation (i.e., at least one citation for the target paper not of the form [X,a,b,c]);
2. The oldest and newest citations within AAN; and,
3. Citations from different years.

We thus provided the oldest and newest citation regardless of criteria 1) and 3) and included a randomized sample of up to 8 additional citing paper

<sup>2</sup><http://www.google.com.sg>

<sup>3</sup><http://www.scholar.google.com>

<sup>4</sup><http://clair.eecs.umich.edu/aan/index.php>

IDs that met either criteria 1) and 3). To do this, we started by first randomizing the list of citing papers and enumerating up to 8 additional citing papers. At this point, the citing papers were listed as either “old(est)”, “new(est)”, 1-8 (additional citing papers), or “sub” (substitute backup citing paper in case of disqualification of one of the 1-8 additional papers due to criterion 1). We combed through the lists of the additional 1-8 citing papers per target paper, from the top to the bottom of the randomized list. We unilaterally collected the top most oldest and top most newest paper, in case of ties. For the remaining (up to) 8 papers, we examined the citing paper’s PDF file to ensure that at least one citation made was of a single citation format (e.g., [X] and not [X,a,b,c]). Any invalidated files were marked with “list” (citation) mark. The resulting final list was divided among the annotators to add human annotations using the same scheme used by annotators of the BiomedSumm track’s corpus.

Given a reference paper and up to 10 citing papers, annotators from National University of Singapore and Nanyang Technological University were instructed to find citations to the reference paper (RP) in the Citing Papers (CP). Annotators followed instructions used for annotation of corpus for the BiomedSumm to encourage cross participation across the two tasks. Specifically, the citation text, citation marker, reference text, and discourse facet were marked for each citation of the RP found in the CP.

## 4 The Task

This Shared Task proposes to solve the problems posed in the BioMedSumm track, but in the domain of Computational Linguistics. This task calls for summarization frameworks to build a structured summary of a research paper - which incorporates faceted information (such as Aims, Methods, Results and Implications) from the text of the paper, and “community summaries” from its citing papers.

The CL-Summ Task is defined as follows:

Given: Ten topics, which comprise a Reference Paper (RP) and up to ten papers which cite it (Citing Papers, or CPs). In every CP, the citations to the RP (known as “citances”) have been identified. The information referenced in the RP has also been iden-

tified in the hand-annotated gold standard version.

Task 1a: Develop a method to identify the text span in the RP which corresponds to the citances from the CP. These may be of the granularity of a full sentence or several sentences, whether contiguous or non-contiguous. It may also be a sentence fragment (no more than 5).

Task 1b: Develop a method to identify the facet for every cited text span from a predefined set of facets.

Evaluation: Evaluate Task 1 by using the ROUGE score to compare the overlap of text spans in the system output vs the gold standard created by human annotators.

## 5 Participating teams

Nine teams expressed an interest in participating, of which two teams have submitted their findings thus far:

1. The MQ System, from Macquarie University, Australia<sup>5</sup>. This system is the same one that was used for the BiomedSumm track, with the exception that it did not incorporate domain knowledge (UMLS). For task 1a it used similarity metrics to extract the top n sentences from the documents. For task 2, they incorporated the distances from task 1 to rank the sentences. Details of their evaluation results are provided in this paper.
2. clair\_umich from University of Michigan, Ann Arbor, USA. This is a supervised system which used lexical and syntactic dependencies as features. Their results are discussed in this paper

Other teams to have expressed an interest are:

1. Taln.UPF, from Universitat Pompeu Fabra, Spain<sup>6</sup>. In the current version of this paper, we have described the algorithm of their approach - they aim to share their results in the

<sup>5</sup>This research was made possible thanks to a summer internship granted to Christopher Jones by the Department of Computing, Macquarie University

<sup>6</sup>This research is supported by the project Dr. Inventor (FP7-ICT-2013.8.1 611383), programa Ramón y Cajal 2009 (RYC-2009-04291), and the project TIN2012-38584-C06-03 Ministerio de Economía y Competitividad, Secretaría de Estado de Investigación, Desarrollo e Innovación, Spain.

near future. They have proposed to adapt available summarisation tools to scientific texts.

2. TXSUMM, from University of Houston, Texas. Their system consists of applying similarity kernels in an attempt to better discriminate between candidate text spans (with sentence granularity). They are using an extractive procedure with ranking algorithms.
3. IITKGP\_sum, from Indian Institute of Technology, Kharagpur, India. They plan to use citation network structure and citation context analysis to summarise the scientific articles.
4. CCS2014, from the IDA Center for Computing Sciences, USA. They will employ a language model based on the sections of the document to find referring text and related sentences in the cited document.
5. TabiBoun14, from the Boazii University, Turkey. They plan to modify an existing system for CL papers, wherein they use LIBSVM as a classification tool for face classification. They also plan to use the cosine similarity metric to compare text spans.
6. PolyAF, from The Hong Kong Polytechnic University.
7. The IHMC system, from IHMC, USA.

## 6 The MQ System - Finding the Best Fit to a Citance

Given the text of a citance, the MQ system ranks the sentences of the reference paper according to its similarity to the citance. Every sentence and its citance was modeled as a vector and compared using cosine similarity. The team experimented with different forms of representing the information in the vectors, and different forms of using the similarity scores to perform the final sentence ranking.

### 6.1 Baseline - Using *tf.idf*

For the baseline system, the *tf.idf* of all lowercased words was used, without removing stop words. Separate *tf.idf* statistics were computed for each reference paper, using the set of sentences in the paper and the citance text of all citing papers.

### 6.2 Adding texts of the same topic

Since the amount of text used to compute the *tf.idf* in Section 6.1 was relatively little, the complete text of all citing papers was added, under the presumption that citing papers are presumably of the same topic as the reference paper. By adding this text we hope to include complementary information that can be useful for extending and computing the *idf* component.

### 6.3 Adding context

In order to extend the information of each sentence in the reference paper and further add to the approach in Section 6.2, the text from the reference papers was added within a context window of 20 sentences by including the neighbouring sentences, centered in the target sentence.

### 6.4 Re-ranking using MMR

The last experiment used Maximal Marginal Relevance (MMR) (Carbonell and Goldstein, 1998) to rank the sentences. All sentences were represented as *tf.idf* vectors of extended information as described in Section 6.3. Then, the final score of a sentence was the combination of the similarity with the citance and similarity of the other sentences of the summary according to the formula shown in Figure 1. A value of  $\lambda = 0.97$  was chosen.

## 7 The clair\_umich System - Comparing Overlap of Word Synsets

### 7.1 Data Preprocessing

The original SciSumm corpus contained data for 10 papers sampled from the ACL Anthology. For each of these papers, citing sentences were extracted from all its citing papers. Each citing sentence was then matched to a text segment in the original paper creating the final annotated dataset. The original source text for the papers in the SciSumm corpus was not sentence segmented, which made it difficult to compute evaluation metrics.

Data preprocessing of the SciSumm corpus was done in the following way - First, sentences from the reference papers were segmented and then matched to each of these source sentences to the SciSumm annotation files. This yielded a fixed set of source sentences from the original files, a subset of which

$$\text{MMR} = \arg \max_{D_i \in R \setminus S} \left[ \lambda(\text{sim}(D_i, Q)) - (1 - \lambda) \max_{D_j \in S} \text{sim}(D_i, D_j) \right]$$

Where:

- $Q$  is the citance text.
- $R$  is the set of sentences in the document.
- $S$  is the set of sentences that haven been chosen in the summary so far.

Figure 1: Maximal Marginal Relevance (MMR)

were matched to each citing sentence. In this way, given a citing sentence, matching sentences from the source paper were compared to the gold standard sentences matched from the source paper and compute precision / recall.

The average number of source sentences matched for each citing sentence was 1.28 (with standard deviation 1.92). The maximum number of source sentences matched for a citing sentence was 7. Given that the total number of source sentences for papers ranged from between 100 to 600, this made it a very challenging classification problem.

## 7.2 Baseline System

Like the MQ system, the team first created a baseline system based on TF\*IDF cosine similarity. For any citing sentence, the system computed the TF\*IDF cosine similarity with all the sentences in the reference paper, thus the IDF values differed across each of the 10 reference papers.

## 7.3 Supervised System

The supervised system used knowledge based features derived from WordNet, syntactic dependency based features, and distributional features in addition to the simple lexical features like cosine similarity. These features are described below.

**Lexical Features** Two lexical features were used - tf\*idf and the LCS (Longest Common Subsequence) between the citing sentence ( $C$ ) and source sentence  $S$ , which is computed as:

$$\frac{|LCS|}{\min(|C|, |S|)}$$

**Knowledge Based Features** The system also used set of features based on Wordnet similarity. Six wordnet based word similarity measures were combined to obtain six knowledge based sentence similarity features using the method proposed in (Banea et al., 2012). The wordnet based word similarity measures used are path similarity, WUP similarity (Wu and Palmer, 1994), LCH similarity (Leacock and Chodorow, 1998), Resnik similarity (Resnik, 1995), Jiang-Conrath similarity (Jiang and Conrath, 1997), and Lin similarity (Lin, 1998).

Given each of these similarity measures, the similarities between two sentences was computed by first creating a set of senses for each of the words in each of the sentences. Given these two sets of senses, the similarity score between citing sentence  $C$  and source sentence  $S$  was calculated as follows:

$$\text{sim}_{wn}(C, S) = \frac{(\omega + \sum_{i=1}^{|\phi|} \phi_i) * (2|C||S|)}{|C| + |S|}$$

Here  $\omega$  is the number of shared senses between  $C$  and  $S$ . The list  $\phi$  contains the similarities of non-shared words in the shorter text,  $\phi_i$  is the highest similarity score of the  $i$ th word among all the words of the lower text (S13, 2013).

**Syntactic Features** An additional feature based on similarity of dependency structures was used, by applying the method described in (S13, 2013). The Stanford parser was used to obtain dependency parse all the citing sentences and source sentences. Given a candidate sentence pair, two syntactic dependencies were considered equal if they have the same dependency type, governing lemma, and dependent lemma. If  $R_c$  and  $R_s$  are the set of all dependency

relations in  $C$  and  $S$ , the dependency overlap score was computed using the formula:

$$sim_{dep}(C, S) = \frac{2 * |R_c \cap R_s| * |R_c| |R_s|}{|R_c| + |R_s|}$$

## 8 The TALN.UPF System

This section details the algorithm to be used in the TALN.UPF system. The results will be included in a future version of this paper.

### 8.1 Pre-processing / documents preparation:

The TALN.UPF system carried out the following set of preprocessing steps on the papers of each topic:

- Sentence segmentation: To identify candidate sentences that will be validated or rejected in the following pre-processing steps;
- Tokenizer and POS tagger: Using the open-source GATE software
- Sentence sanitizer: To remove incorrectly annotated sentences, relying on a set of rules and heuristics;
- Document structural analyzer: To classify each sentence as belonging to one of the following document structural categories: Abstract, Introduction, Result\_Discussion, Experimental\_Procedure, Supplemental\_Data, Material\_Method, Conclusion, Acknowledgement\_Funding, and Reference;
- Sentence TFIDF vector computation: To associate to each sentence a TFIDF vector where the IDF values are computed over all the papers of the related topic (up to 10 citing paper and one reference paper).
- From the citing paper, those sentences were selected which overlapped totally or partially the global citance context span; these sentences were referred to as the citance context sentences (CtxSent1,..., CtxSentN),
- Citances were characterized by the document structural category associated with most of its citance context sentences (CtxSent1,..., CtxSentN). In case of tie in the number of occurrences of document structural categories among all the citance context sentences, the most frequently chosen document structural category for the citing paper was preferred. In case of persisting ties, the document structural category that is most frequent in the whole set of citing and reference papers was preferred.
- Each reference paper sentence (RefSent) was assigned a score equal to the sum of its TF\*IDF vector cosine similarity with each citance context sentence (CtxSent1,..., CtxSentN).
- The RefSent scores were weighted by the relative relevance of this kind of link between document structural categories, in the whole training corpus. For instance, if there is a citance associated to the INTRO that references a RefSent belonging to the Abstract and in the whole training corpus this situation occurs in 6.5% of citance-referenced sentence pairs, the RefSent score is multiplied by 0.065, obtaining the final RefSent score.
- The first 3 reference paper sentences (RefSents) with the highest final RefSent score were chosen as the reference paper text spans.

### 8.2 Algorithm for identifying reference paper text spans for each citance

- For each citance its global citance context span was considered as the union of the citance context spans marked by human annotators (in this case, there was only one available human annotation, so no union was required).

### 8.3 Algorithm for identifying the discourse facet of the cited text spans

A linear-kernel SVM classifier was trained to associate each citance with one of the five text facets considered in Task 1b. Each citance was characterized by lexical and semantic features extracted from the sentences belonging to the citance context together with the sentences of the reference paper selected as outcome of Task 1a. Some of the features exploited were:

1. Relative number of sentences belonging to each document structural category
2. Relative number of sentences belonging to the citance context or reference paper
3. Relative number of POS
4. Presence of key lexical patterns

## 9 Evaluation and Results

Two of the teams have submitted their results so far, and the evaluation is based on the ROUGE metric (Lin, 2004). ROUGE is a popular evaluation method for summarisation systems that compares the text output of the system against a set of target summaries. Since ROUGE uses the actual contents words, and not the offset information, we expect that this metric will give non-zero results for cases when a system chooses a sentence that is similar to, but not exactly, the one chosen by the annotator.

The MQ system and clair\_umich system were both unsupervised, so for the evaluation, they were able to use all the data without having to perform cross-validation experiments. For the MQ system, the output is the set of selected sentences, and the target summaries are the sentences given by the annotators. For the clair\_umich system, the ROUGE-L scores were computed for each citing sentence in each annotation file separately and then averaged for a topic.

The following paragraphs describe the results for Task 1a, 1b, and the bonus Task 2 which was attempted by the MQ system.

### 9.1 Task 1a: For each citance, identify the spans of text (cited text spans) in the RP

Table 2 shows the ROUGE-L F1 scores of each individual reference document from the SciSumm dataset.

### 9.2 Task 2: Generate a structured summary of the RP and all of the community discussion of the paper represented in the citances

The MQ team performed an additional test to see whether information from the citances were useful for building an extractive summary, as is the case with the BiomedSumm data (Mollá et al., 2014).

They implemented extractive summarisation systems with and without information from the citances. The summarisers without information from the citances scored each sentence as the sum of the *tf.idf* values of the sentence elements. They tried the *tf.idf* approach described in Section refsec:tfidf.

The summarisers with information from the citances scored each candidate sentence  $i$  on the basis of  $\text{rank}(i, c)$  obtained in task 1a, which has values between 0 (first sentence) and  $n$  (last sentence) and represents the rank of sentence  $i$  in citance  $c$ :

$$\text{score}(i) = \sum_{c \in \text{citances}} 1 - \frac{\text{rank}(i, c)}{n}$$

The summaries were evaluated using ROUGE-L, where the model summaries are the abstract section of the corresponding papers. Since paper X96-1048 of the SciSumm data did not have an abstract section, it was removed for this experiment.

An example excerpt from a target summary (Abstract) for the reference paper J03-3003 is:

*We describe a statistical approach for modeling dialogue acts in conversational speech, i.e., speech-act-like units such as STATEMENT, QUESTION, BACKCHANNEL, AGREEMENT, DISAGREEMENT, and APOLOGY. Our model detects and predicts dialogue acts based on lexical, collocational, and prosodic cues, as well as on the discourse coherence of the dialogue act sequence. The dialogue model is based on treating the discourse structure of a conversation as a hidden Markov model and the individual dialogue acts as observations emanating from the model states. Constraints on the likely sequence of dialogue acts are modeled via a dialogue act n-gram... We achieved good dialogue act labeling accuracy (65% based on errorful, automatically recognized words and prosody, and 71% based on word transcripts, compared to a chance baseline accuracy of 35% and human accuracy of 84%) and a small reduction in word recognition error.*

The MQ System's output baseline summary for the same reference paper is 20 sentences along; below is an excerpt: *Dialogue Act Modeling for Automatic Tagging and Recognition of Conversational Speech. In all these cases, DA labels would enrich the available input for higher-level processing*

	MQ System			clair_umich System		
Run	R	P	F1	R	P	F1
Using all features	0.335	0.212	0.223	0.0	0.0	0.59

Table 1: ROUGE-L results of the participating systems for task 1a

Paper ID	MQ System	clair_umich System
C90-2039	0.235	0.698
C94-2154	0.288	0.638
E03-1020	0.239	0.579
H05-1115	0.350	0.697
H89-2014	0.332	0.658
J00-3003	0.196	0.484
J98-2005	0.101	0.656
N01-1011	0.221	0.603
P98-1081	0.200	0.531
X96-1048	0.248	0.410

Table 2: ROUGE-L F1 results for individual topics 1a

of the spoken words. The relation between utterances and speaker turns is not one-to-one: a single turn can contain multiple utterances, and utterances can span more than one turn (e.g., in the case of backchanneling by the other speaker in midutterance). The most common of these are the AGREEMENT/ACCEPTS. One frequent example in our corpus was the distinction between BACKCHANNELS and AGREEMENTS (see Table 2), which share terms such as “right” and “yeah”. Networks compare to decision trees for the type of data studied here. Neural networks are worth investigating since they offer potential advantages over decision trees.

Table 3 shows the breakout of ROUGE-L F1 scores per document.

## 10 Discussion

### 10.1 Comparing the MQ System with the BioMedSumm task

Table 4 compares the results of the MQ system’s experiments with the SciSumm data, against the results from the BiomedSumm data. In all results the systems were designed to return 3 sentences, as specified in the shared task. All short sentences (under 50 characters) were ignored, to avoid including headings or mistakes made by the sentence segmentation algorithm.

The results show an improvement in both domains, with the exception that MMR does not improve over the run that uses *tf.idf* over context in SciSumm, whereas there is an improvement in BiomedSumm. The absolute values are better in the BiomedSumm data, and looking at the confidence intervals it can be presumed that the difference between the best and the worst run is statistically significant in the BiomedSumm data. The results in the SciSumm data are poorer in general and there are no statistically significant differences. However, this may be an artifact of the small size of the corpus. Overall, the improvement of results in SciSumm mirrors that of the BiomedSumm data, so it can be suggested that on adding more information to the models that compute *tf.idf*, the results improve. It is expected that alternative approaches, which gather related information to be added for computing the vector models will produce even better results. The results with MMR appears to be contradictory across the two domains but the difference is so small that it might not be statistically significant even when we add more evaluation data.

These experiments suggest that information from the citances may be useful for building an extractive summary. This conclusion is compatible with prior research that suggest that, in general, informa-



Paper ID	<i>tf.idf</i>	task1a <i>tf.idf</i>	task1a MMR
C90-2039_TRAIN	0.347	0.315	0.293
C94-2154_TRAIN	0.095	0.123	0.120
E03-1020_TRAIN	0.189	0.189	0.196
H05-1115_TRAIN	0.134	0.306	0.321
H89-2014_TRAIN	0.294	0.319	0.320
J00-3003_TRAIN	0.221	0.382	0.367
J98-2005_TRAIN	0.221	0.216	0.233
N01-1011_TRAIN	0.187	0.268	0.284
P98-1081_TRAIN	0.241	0.210	0.206
Average	0.214	0.259	0.260

Table 3: ROUGE-L F1 results for summaries generated by the MQ system

Run	SciSumm				BiomedSumm			
	R	P	F1	CI	R	P	F1	CI
<i>tf.idf</i>	0.316	0.198	0.211	0.185–0.240	0.273	0.326	0.279	0.265–0.293
topics	0.324	0.201	0.217	0.191–0.245	0.288	0.357	0.300	0.285–0.316
context	0.339	0.214	0.225	0.197–0.255	0.291	0.372	0.308	0.293–0.323
MMR	0.335	0.212	0.223	0.195–0.251	0.290	0.375	0.308	0.293–0.323

Table 4: ROUGE-L results of the MQ system runs for task 1a

tion from citing papers may be useful for building summaries, as was stated in the original goals of the Sci-Summ Shared Task.

## 10.2 Error Analysis for the MQ System

Some drawbacks were observed in the approach and evaluation for the MQ system. The example below illustrates the MQ system’s output for task1a, for the reference paper H89-2014:

“The statistical methods can be described in terms of Markov models.” “An alternative approach taken by Jelinek, (Jelinek, 1985) is to view the training problem in terms of a “hidden” Markov model: that is, only the words of the training text are available, their corresponding categories are not known.” “In this regard, word equivalence classes were used (Kupiec, 1989).”

The target sentence was: “The work described here also makes use of a hidden Markov model.”

The first sentence of the sample output was very similar to the target sentence. It was not the best match, but it was a close match, and an evaluation metric such as ROUGE would reward it. On the other hand, the second sentence, even though it

talked about HMMs, it was not strictly about the approach used by the paper and therefore it should not be rewarded with a good score. However, ROUGE would be too lenient here. This is one of the issues identified by the MQ system in following a purely lexical approach.

## 10.3 Tweaking the Parameters - the clair\_umich Baseline

For any citing sentence, the TF\*IDF cosine similarity was computed with all the sentences in the source paper, and any sentences that had a cosine similarity higher than a given threshold were added to the matched sentences. Table 5 shows the precision/recall for different values of the cosine threshold:

The F1-score seems to reach a maxima at a similarity threshold of about 0.1. The recall at the threshold of 0.1 is about 0.23, while the precision is only 0.06. This suggests that initial progress can be made on this problem by first removing these spurious matches that have high lexical similarity.

**Error Analysis for the clair\_umich Baseline System** A number of errors made by the baseline sys-

Similarity Threshold	Precision	Recall	F1-score
0.01	0.027	0.641	0.051
0.05	0.048	0.426	0.087
0.1	0.060	0.235	0.095
0.2	0.079	0.081	0.080
0.3	0.062	0.032	0.042
0.4	0.022	0.085	0.012
0.5	0.007	0.002	0.003

Table 5: Precision/Recall for different values of the cosine threshold for the baseline clair\_umich system

tem are due to source sentences that match the words but differ slightly in their information content. Here is an example.

Citing text: “use the BNC to build a co-occurrence graph for nouns, based on a co-occurrence frequency threshold”

*True positives:*

- “Following the method in (Widdows and Dorow, 2002), we build a graph in which each node represents a noun and two nodes have an edge between them if they co-occur in lists more than a given number of times.”

*False positives:*

- “Based on the intuition that nouns which co-occur in a list are often semantically related, we extract contexts of the form Noun, Noun,... and/or Noun, e.g. “genomic DNA from rat, mouse and dog”.”
- “To detect the different areas of meaning in our local graphs, we use a cluster algorithm for graphs (Markov clustering, MCL) developed by van Dongen (2000).”
- “The algorithm is based on a graph model representing words and relationships between them.”

Even though the false positive sentences contain the same lexical items (nouns, co-occurrence, graph), they differ slightly in the facts presented. Detection of such subtle differences in meaning might be challenging for an automated system.

Another set of difficult sentences is when the citing sentence says something that is implied by the sentence in the source paper. For example:

Citing text: “The line of our argument below follows a proof provided in ... for the maximum likelihood estimator based on nite tree distributions”

*False negatives:*

- “We will show that in both cases the estimated probability is tight.”

Here, the citing text mentions a proof from source paper, but to match the sentence in the source paper, the system needs to understand that the act of showing something in a scientific paper constitutes a proof.

## 11 Limitations and Error Analysis

There were several limitations in the dataset which were identified in the process of annotating and parsing the corpus for use by the participating systems; these are discussed below.

- The use of “...” where text spans are snippets: The use of “...” follows the BioMedSumm standard practice of indicating discontinuous texts. In Citation Text and Reference Text fields, the “...” means that there is a gap between two text spans (citation spans or reference spans). They may be on different pages, so the gap might be a text. There might be a formula or a figure there, or some text encoding which is not a part of the annotation. However, this notation caused mismatches for sentences which used text from different parts of the same sentence.
- Small size of the training corpus: The corpus comprised only a training set of 10 topics, each with upto 10 citing documents. In this small

dataset, participants were asked to conduct a 10-fold cross validation. The small size of the data set meant that there were no statistically significant results, but significance could only be guessed at, from the overall trend of the data.

- Errors in parsing the file: Some of the older PDF files, when parsed to text or xml, had such as misspelt words, spaces within words, sentences in the wrong place and so on. Unfortunately these errors were OCR parsing errors, and not in our control. It was recommended that the participants should configure their string matching to be lenient enough to tackle such problems.
- Errors in citation/reference offset numbers: In the original annotations, citation/reference offset numbers were character-based, and relative to an xml encoding which was not shared in the task, and did not match with the offset numbers on the text-only, cleaned version of the document. Although the text versions of the source documents were shared with the intention to help the participants, this often made their tasks more difficult if their system was geared towards numerical and not system matching. A solution was found for reference offsets by revising them to sentence id numbers based on available XML files from the clair\_umich system's pre-processing stage; however, the citation offsets remain character-based.
- Text encoding: Often, the text was not in UTF-8 format as expected. Some participating teams, like the UPF, solved this by running the universal charset tool provided by Google Code over all the text and annotations in order to determine the right file encoding to use. It was found that some of the files were also in WINDOWS-1252 and GB18030 formats.
- Errors in file construction: An automatic, open-source software was used to map the citation annotations from a software, Protege, to a text file. However, participants identified several errors in the output - especially in cases where there was one-to-many mapping between citations and references. Besides this, several

annotation texts had no annotation id (Citance Number field).

## 12 Conclusion

This paper describes the informal SciSumm Computational Linguistics task for faceted summarization of research papers. Two systems participated in the task-based evaluation and submitted their results. The teams used versions of TF\*IDF as baselines, but the MQ system followed an unsupervised approach which clair\_umich followed a supervised approach. For identifying referenced text spans in reference papers, the best performance was obtained by clair\_umich by following a supervised approach using lexical, syntactic and knowledge-based features to calculate the overlap between sentences in the citation span and the reference paper. Although no system submitted results for the the task involving identifying the discourse facets of reference text, TALN.UPF submitted an algorithm which they aim to implement. Finally, an added experiment by the MQ system sought to compare baseline summaries of reference papers, based on a TF\*IDF calculation, against gold standard summaries, comprising the reference paper's abstracts.

The clair\_umich system incorporated WordNet synsets for expanding and comparing cited text with reference papers, and the use of syntactic features further enriched the calculation of overlap. On the other hand, the MQ system relied exclusively on reading and comparing texts. Furthermore, their system was originally built for the BioMedSumm task - however, they had to discard some domain-specific features for this task. It is possible that the lack of domain knowledge, coupled with OCR-related and PDF parsing errors, affected the performance of their system in the CL domain.

This task is an initiative for encouraging the development of tools and approaches for scientific summarisation. It helped in identifying the existing tools and resources which could be leveraged for this purpose, and also the hindrances which need to be overcome in order to have a systematic and well-coordinated evaluation. However, with the results of only two systems, it is not possible to conjecture at what may be the better methods for summarizing research papers in this domain. The resources from the

SciSumm task, and its corpus, are freely available for interested research groups to experiment with; the corpus is first-of-its-kind summarization corpus for computational linguistics.

The SciSumm organization committee hopes to carry this year's efforts forward as a full-fledged task in the coming year. We plan a systematic annotation of a training set, as well as a test set, the availability of more than one gold standard annotation, and open-sourced tools and resources to support the efforts of participating teams. We invite the community to join us in this endeavour with any resources and time they can spare.

## References

- Amjad Abu-Jbara and Dragomir Radev. 2011. Coherent citation-based summarization of scientific papers. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, pages 500–509. Association for Computational Linguistics.
- Carmen Banea, Samer Hassan, Michael Mohler, and Rada Mihalcea. 2012. Unt: A supervised synergistic approach to semantic text similarity. In *Proceedings of the First Joint Conference on Lexical and Computational Semantics-Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation*, pages 635–642. Association for Computational Linguistics.
- Steven Bird, Robert Dale, Bonnie J Dorr, Bryan R Gibson, Mark Joseph, Min-Yen Kan, Dongwon Lee, Brett Powley, Dragomir R Radev, and Yee Fan Tan. 2008. The acl anthology reference corpus: A reference dataset for bibliographic research in computational linguistics. In *LREC*.
- Jaime Carbonell and Jade Goldstein. 1998. The use of MMR, diversity-based reranking for reordering documents and producing summaries. In *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval - SIGIR '98*, pages 335–336, New York, New York, USA. ACM Press.
- Cong Duy Vu Hoang and Min-Yen Kan. 2010. Towards automated related work summarization. In *Proceedings of the 23rd International Conference on Computational Linguistics: Posters*, pages 427–435. Association for Computational Linguistics.
- Kokil Jaidka, Christopher SG Khoo, and Jin-Cheon Na. 2013. Literature review writing: how information is selected and transformed. In *Aslib Proceedings*, volume 65, pages 303–325. Emerald Group Publishing Limited.
- J.J. Jiang and D.W. Conrath. 1997. Semantic similarity based on corpus statistics and lexical taxonomy. In *Proc. of the Int'l. Conf. on Research in Computational Linguistics*, pages 19–33.
- C. Leacock and M. Chodorow. 1998. Combining local context and wordnet similarity for word sense identification. In Christiane Fellbaum, editor, *MIT Press*, pages 265–283, Cambridge, Massachusetts.
- Dekang Lin. 1998. An information-theoretic definition of similarity. In *Proceedings of the Fifteenth International Conference on Machine Learning, ICML '98*, pages 296–304, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Chin-Yew Lin. 2004. {ROUGE}: A Package for Automatic Evaluation of Summaries. In *ACL Workshop on Tech Summarisation Branches Out*.
- Saif Mohammad, Bonnie Dorr, Melissa Egan, Ahmed Hassan, Pradeep Muthukrishnan, Vahed Qazvinian, Dragomir Radev, and David Zajic. 2009. Using citations to generate surveys of scientific paradigms. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 584–592. Association for Computational Linguistics.
- Diego Mollá, Christopher Jones, and Abeed Sarker. 2014. Impact of Citing Papers for Summarisation of Clinical Documents. In *Proc. ALTA 2014*.
- Philip Resnik. 1995. Using information content to evaluate semantic similarity in a taxonomy. In *Proceedings of the 14th International Joint Conference on Artificial Intelligence - Volume 1, IJCAI'95*, pages 448–453, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
2013. *ECNUS: Measuring Short Text Semantic Equivalence Using Multiple Similarity Measurements*, Second Joint Conference on Lexical and Computational Semantics (\*SEM), Volume 1: Proceedings of the Main Conference and the Shared Task: Semantic Textual Similarity. Association for Computational Linguistics.
- HOANG CONG DUY VU. 2010. Towards automated related work summarization.
- Zhibiao Wu and Martha Palmer. 1994. Verbs semantics and lexical selection. In *Proceedings of the 32Nd Annual Meeting on Association for Computational Linguistics, ACL '94*, pages 133–138, Stroudsburg, PA, USA. Association for Computational Linguistics.