# The use of domain-specific concepts in biomedical text summarization

Lawrence H. Reeve [a], Hyoil Han [a,*], Ari D. Brooks [b]

[a] *Drexel University, College of Information Science and Technology, Philadelphia, PA, USA*
[b] *Drexel University, College of Medicine, Philadelphia, PA, USA*

## Abstract

Text summarization is a method for data reduction. The use of text summarization enables users to reduce the amount of text that must be read while still assimilating the core information. The data reduction offered by text summarization is particularly useful in the biomedical domain, where physicians must continuously find clinical trial study information to incorporate into their patient treatment efforts. Such efforts are often hampered by the high-volume of publications. This paper presents two independent methods (*BioChain* and *FreqDist*) for identifying salient sentences in biomedical texts using concepts derived from domain-specific resources. Our semantic-based method (BioChain) is effective at identifying thematic sentences, while our frequency-distribution method (FreqDist) removes information redundancy. The two methods are then combined to form a hybrid method (*ChainFreq*). An evaluation of each method is performed using the ROUGE system to compare system-generated summaries against a set of manually-generated summaries. The BioChain and FreqDist methods outperform some common summarization systems, while the ChainFreq method improves upon the base approaches. Our work shows that the best performance is achieved when the two methods are combined. The paper also presents a brief physician's evaluation of three randomly-selected papers from an evaluation corpus to show that the author's abstract does not always reflect the entire contents of the full-text.
© 2007 Published by Elsevier Ltd.

*Keywords:* Text summarization; Biomedicine; Concept chaining; Concept frequency

## 1. Introduction

Text summarization is a data reduction process. The use of text summarization allows a user to get a sense of the content of a full-text, or to know its information content, without reading all sentences within the full-text. Data reduction increases scale by (1) allowing users to find relevant full-text sources more quickly, and (2) assimilating only essential information from many texts with reduced effort. Text summarization is particularly useful in the biomedical domain, where oncologists must continuously find clinical trial

32  study information related to their specialty, evaluate the study for its strength, and then possibly incorporate
33  the new study information into their patient treatment efforts (Brooks & Sulimanoff, 2002; Jaques, 2002). The
34  US National Institutes of Health Clinical Trials database contains over 13,500 clinical trials (United States
35  National Library of Medicine, 2005a), while PUBMED contains in excess of 16 million citations from over
36  4800 journals (United States National Library of Medicine, 2006a). Large and continuously-updated informa-
37  tion sources such as these impede a physician's ability to improve their treatment efforts.
38      The contribution of this work is to present three novel extractive summarization methods (BioChain for
39  concept chaining, FreqDist for frequency distribution, and a hybrid of the two, ChainFreq) using domain-spe-
40  cific concepts as a feature for identifying salient sentences in biomedical texts, and present an evaluation of
41  them along with several other publicly-available summarizers. We show that our semantic-based summarizer
42  (BioChainSumm) is effective at identifying thematic sentences, while our frequency-distribution summarizer
43  (FreqDistSumm) removes information redundancy. The best performance is achieved when the two methods,
44  BioChain and FreqDist are combined to form a third hybrid summarizer, ChainFreqSumm.
45      The paper is organized as follows. Section 2 provides related work and background on the need for bio-
46  medical text summarization and the methods used for concept annotation. Section 3 presents our three con-
47  cept-based algorithms. Section 4 describes the evaluation methodology and Section 5 discusses the results of
48  the evaluation and presents a physician's observations on the mismatch between abstracts and full-text of sev-
49  eral selected papers. Section 6 provides concluding remarks and identifies future work.

50  **2. Background and related work**

51  *2.1. Need for biomedical text summarization*

52      Clinical trial studies and other scientific publications usually supply a summary of the paper in the form of
53  an abstract produced by the author(s) of a study. We have identified at least five reasons for wanting to gen-
54  erate text summaries from a full-text source even in the presence of the author's abstract. (1) There exists no
55  'ideal' summary. An ideal summary is dependent on each user, including factors such as information need and
56  domain background. An author's abstract is one view of an ideal summary, but users may want alternative
57  summaries. (2) The abstract may be missing content from the full-text (Cohen & Hersh, 2005). (3) Customized
58  summaries can be useful in question-answering systems where they provide personalized information. (4) The
59  use of automatic or semi-automatic summary generation by commercial abstract services may allow them to
60  scale the number of published texts they can evaluate. (5) The generation and evaluation of summaries allows
61  for evaluation of sentence selection methods that may be useful in multi-document summarization.

62  *2.2. Biomedical domain concepts and automated concept annotation*

63      The main idea of this work is to use domain-specific (biomedical) concepts to identify important sentences
64  within a biomedical text which can be extracted to form a summary. To achieve text summarization through
65  the use of concepts, three resources are required: (a) a list of domain-specific concepts, (b) one or more syn-
66  onymous phrases which occur in text and are associated with each domain-specific concept, and (c) an auto-
67  mated method for identifying concepts within a text. Concept annotation of each paper in the evaluation
68  corpus was performed using the UMLS MetaMap Transfer application (United States National Library of
69  Medicine, 2005b). Summary generation using the discovered concepts then takes place in two stages: (1) bio-
70  medical concept annotation of the source text, and (2) summary generation from the concept-annotated text
71  using the discovered concepts.
72      In the biomedical domain, the National Library of Medicine (http://www.nlm.nih.gov/) provides resources
73  for identifying concepts and their relationships under the framework of the Unified Medical Language System
74  (UMLS) (United States National Library of Medicine, 2005c). UMLS contains many sub-components, but we
75  use only three: Metathesaurus, Semantic Network, and MetaMap Transfer. The UMLS Metathesaurus,
76  derived from over 100 different biomedical vocabulary sources, contains concepts and real-world instances
77  of the concepts, including a concept name and its synonyms (United States National Library of Medicine,

78  2006b). The UMLS Semantic Network organizes the Metathesaurus concepts into semantic types (United
79  States National Library of Medicine, 2004).

80      For automated concept annotation, the MetaMap Transfer (MMTx) (United States National Library of
81  Medicine, 2005b) application maps free-form biomedical text to Metathesaurus concepts. MMTx performs
82  text-to-concept mapping by first identifying noun phrases in each sentence, generating term variants of the
83  phrase, finding candidate concepts from the generated phrase variants, and scoring each candidate concept.
84  The highest scoring concept and its semantic type are then returned. It is possible for a noun phrase to
85  map to more than one concept. In this case, no disambiguation is performed, and MMTx returns multiple
86  concepts and their semantic types.

## 2.3. Related work

88      We mainly investigate previous text summarization work related to lexical chaining and frequency meth-
89  odologies because our work largely involves these two approaches.

90      Lexical chaining is a method for determining lexical cohesion among terms in a text (Barzilay & Elhadad,
91  1997), and has been used for many years for text summarization. Lexical cohesion is a property of text that
92  causes a discourse segment to ''hang together'' as a unit (Morris & Hirst, 1991). Lexical cohesion is important
93  in computational text understanding for two major reasons: (1) providing term ambiguity resolution, and (2)
94  providing information for determining the meaning of text (Morris & Hirst, 1991). Lexical chaining is useful
95  for determining the ''aboutness'' of a discourse segment, without fully understanding the discourse. A basic
96  assumption is the text must explicitly contain semantically-related terms identifying the main concept. Lexical
97  chains for text summarization were first introduced by Morris and Hirst (1991). Their initial work described
98  the approach, but did not implement it because electronic versions of a thesaurus were not available at the
99  time. A thesaurus is used to relate words semantically; for example, through synonymy and hypernym/hyp-
100 onym relationships. A machine implementation by Barzilay and Elhadad (1997) showed the theoretical work
101 by Morris/Hirst could be practically realized for document summarization. While Barzilay/Elhadad proved
102 the feasibility of computing lexical chains, their algorithm runs in exponential time. A linear time algorithm
103 was later defined and implemented by Silber and McCoy (2002). A more recent implementation focuses on
104 improving word sense disambiguation based on the idea of one sense per discourse (Galley & McKeown,
105 2003). All of these implementations use WordNet (Fellbaum, 1998) as the knowledge source for identifying
106 semantic relationships between terms. A computational model for semantic relationships between terms
107 was developed by Fellbaum (1998). To our knowledge, no biomedical text summarization used lexical chain-
108 ing with UMLS.

109     Term frequency was first used in extractive text summarization in the late 1950s (Luhn, 1958). A follow-up
110 study of an analysis of five term frequency methods showed high agreement in sentence selection among the
111 methods (Rath, Resnick, & Savage, 1961). Subsequent research using frequency methods focused on the use of
112 frequency as one feature among many for identifying important sentences, such as cue phrases (Edmundson,
113 1999; Pollock & Zamora, 1975). Summarization using larger units of text has also been researched. The LAKE
114 system uses keyphrases for summarization (D'Avanzo, Magnini, & Vallin, 2004). The SUMMARIST system
115 (Hovy & Lin, 1999) uses WordNet (Fellbaum, 1998) concept counting not for identifying salient sentences, but
116 for generalizing concepts for topic interpretation (e.g., {pear, apple} → fruit). Most recently, the SumBasic
117 algorithm uses term frequency as part of a context-sensitive approach to identifying important sentences while
118 reducing information redundancy (Nenkova & Vanderwende, 2005). The use of frequency as a feature in
119 locating important areas of a text has been proven useful in the literature (Edmundson, 1999; Luhn, 1958;
120 Pollock & Zamora, 1975; Rath et al., 1961). This is most likely due to reiteration, where authors state impor-
121 tant information in several different ways, in order to reinforce main points (Sparck Jones, 1999).

## 3. Summarization methods

123     In extractive text summarization, the task is to identify sentences in a source text which are relevant to the
124 user while simultaneously reducing information redundancy. Sentences are scored based on a set of features.

125 The top-*n* highest scoring sentences in a text are then extracted, using *n* as an upper bound, and presented to
126 the user in their order of appearance in the original source text.

## 3.1. Concept chaining (BioChain)

128     The BioChain method (Reeve, Han, & Brooks, 2006) applies the concepts and methods of lexical chaining
129 to biomedical text using concepts rather than terms. Previous lexical chaining approaches use linkages among
130 word instances to identify semantically-related terms, and the resulting linkages are used to identify the themes
131 of a text. In the BioChain method, automatically identified concepts in the source text (see Section 2.2) are
132 chained based on their UMLS semantic type(s). Each concept chain contains a list of concepts belonging
133 to a particular UMLS semantic type. If a concept belongs to multiple UMLS semantic types (i.e., multiple
134 concept chains), the concept appears in multiple concept chains. Fig. 1 shows the concept-chaining summari-
135 zation process. The source text is first processed by the UMLS Metamap Transfer application to identify noun
136 phrases, which are then automatically mapped into UMLS concepts as well as UMLS semantic types. The
137 summarizer (BioChainSumm) then takes all discovered concepts and constructs chains of concepts which
138 are related by their UMLS semantic type.
139     Once all concepts in the text have been chained, the strongest concept chains are identified through scoring.
140 The original lexical chaining paper defines three strong chain features: reiteration, density, and length (Morris
141 & Hirst, 1991). Reiteration is repetition of concepts throughout a text. Density suggests concepts closer
142 together are more likely to be related. Length is the number of concept instances within a concept chain.
143 We score concept chains by multiplying the frequency of the most frequent concept in the concept chain
144 by the number of distinct concepts in the concept chain (Barzilay & Elhadad, 1997; Doran, Stokes, Dunnion,
145 & Carthy, 2004). Once all concept chains are scored, the strongest concept chains are identified. Lexical chain-
146 ing uses two standard deviations above the mean of all chain scores to identify strong chains (Barzilay & Elha-
147 dad, 1997), and we follow that method. Strong concepts within each strong concept chain are then identified
148 using two different methods: (1) using the most frequent concept within each concept chain (ties result in mul-
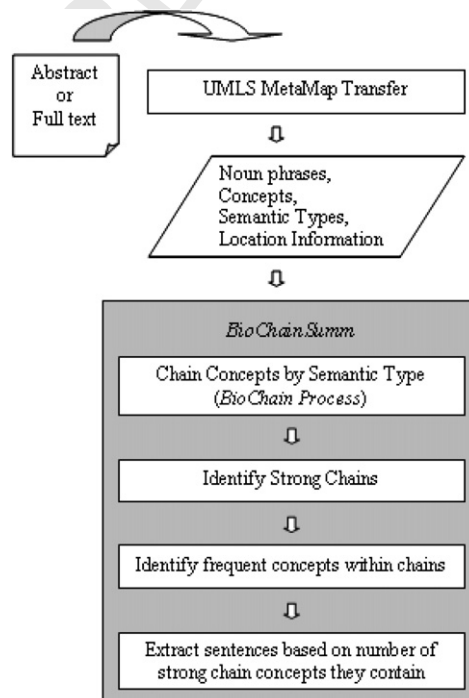


Fig. 1. BioChainSumm summarization process.

```
T081 - Quantitative Concept
        Noun Phrase: high dose intensities
                UMLS Concept: High dose (qualifier value), sentence#3, section#1
        Noun Phrase: primary diameter cm
                UMLS Concept: cm (qualifier value),  sentence#5, section#2
        Noun Phrase: size recurrent tumor
                UMLS Concept: Size (attribute), sentence#5, section#2
        Noun Phrase: One hundred four patients
                UMLS Concept: One (qualifier value), sentence#6, section#2
        Noun Phrase: median disease-free survival
                UMLS Concept: Disease-Free Survival, sentence#9, section#3
        Noun Phrase: median disease-free survival
                UMLS Concept: Median Statistical Measurement, sentence#9, section#3
        Noun Phrase: median overall survival
                UMLS Concept: survival aspects, sentence#10, section#3
        Noun Phrase: median overall survival
                UMLS Concept: Median Statistical Measurement, sentence#10, section#3
        Noun Phrase: absolute benefit
                UMLS Concept: benefits, sentence#11, section#3
```

Fig. 2. A sample concept chain for UMLS Semantic Type #81, quantitative concept.

149 tiple strong concepts for a concept chain) (called MostFrequentStrongChainConcept), and (2) using all con-
150 cepts within a concept chain (called AllStrongChainConcepts). Sentences are scored based on the number of
151 strong concepts they contain. A subset of the original source text sentences is then extracted based on their
152 score and then presented in their original presentation order to form a summary. The size of the generated
153 summary is user-controlled.
154     Fig. 2 shows an example concept chain which links together concepts found within a biomedical source
155 text. In this example, the concept chain is for UMLS semantic type #81, which has a description of *Quanti-*
156 *tative Concept*. The list of concept chain member phrases are phrases automatically found by the MetaMap
157 Transfer application. Each phrase also shows the UMLS concept which maps to it, as well as the sentence
158 and section (paragraph) in which it occurs. This information is also automatically generated by the MetaMap
159 Transfer application.

160 *3.2. Frequency distribution (FreqDist)*

161     When using frequency as the only feature for identifying salient sentences, unit items (e.g., word, concept,
162 and phrase) are counted and then each sentence is given a score based on the frequency count of each unit item
163 in the sentence. A key issue in generating summaries is reducing redundancy. Each new sentence in the sum-
164 mary should add new information. Using the highest frequency terms will likely result in the same information
165 repeatedly being selected. In the SumBasic context sensitive approach (Nenkova & Vanderwende, 2005), a
166 term probability is first determined, and then each term's probability is reduced as it is used to select sentences.
167 This is also related to the idea of Maximal Marginal Relevance (MMR), where marginal relevance is finding
168 relevant sentences which contain minimal similarity to previously selected sentences (Carbonell & Goldstein,
169 1998). We use a context sensitive approach to scoring sentences based on a frequency distribution model
170 rather than term probability. The rationale of our approach is that the frequency distribution of terms or con-
171 cepts in the source text and the generated summary should be as similar as possible.
172     Our FreqDist method (Reeve, Han, Nagori et al., 2006) uses a frequency distribution approach with two
173 stages: Initialization and Summary Generation. In the Initialization stage, the unit items (terms, concepts, etc.)
174 of the source text are counted to form a frequency distribution model of the text, and a pool of sentences from
175 the source text is created. A summary frequency distribution model is then created from the unit items found
176 in the source text, and their frequency counts are initialized to zero. In the Summary Generation stage, new
177 sentences are selected to be added to the summary. Identifying the next sentence to be added to the summary is
178 accomplished by finding the sentence which most closely aligns the frequency distribution of the summary to

```
Initialization:
// Note: '-model' means  'frequency distribution model'
INITIALIZE source-model to unit-items in source-text;
INITIALIZE summary-model,
                candidate-model from source-model;
             set all frequency values of both models to 0;
INITIALIZE sentence-pool to source-text sentences;

Summary Generation:
REPEAT
  INITIALIZE sentence-pool scores to 0;
  INITIALIZE best-score to 0;
  INITIALIZE best-sentence to first sentence in pool;
  INITIALIZE summary-output to empty sentence list;

  FOR each sentence-entry in sentence-pool
    INITIALIZE candidate-model from summary-model;
    ADD sentence unit-item frequencies to candidate-model;
    SET sentence-entry.score =
            similarity(source-model, candidate-model);

   IF sentence-entry.score > best-score
     SET best-score to sentence-entry.score;
     SET best-sentence to sentence-entry;
   ENDIF
  ENDFOR

  ADD unit-items from best- sentence to summary-model;
  ADD best-sentence to summary-output;
  REMOVE best-sentence from sentence-pool;
UNTIL desired summary size reached or
         sentence-pool exhausted;
RETURN summary-output as a final summary;
```

Fig. 3. FreqDistSumm summarization algorithm.

the frequency distribution of the original source text. For each sentence in the sentence pool, a candidate summary is first initialized to the summary generated so far, and then the sentence is added to the candidate summary. The candidate summary frequency distribution is then compared for similarity to the original source text frequency distribution. This similarity score (see next paragraph) is assigned to the sentence. That is, the similarity calculation step is applied to each sentence in the sentence pool one-by-one. After all sentences from the sentence pool have been evaluated for their contribution to the candidate summary, the highest scoring sentence is added to the summary and removed from the sentence pool. This process of identifying the next sentence to be added to the summary is iterative, and repeats until the desired length of the summary is reached. Fig. 3 shows the complete FreqDist method implemented in the FreqDistSumm summarizer.

We compared five similarity functions to find which type of function worked best to evaluate a candidate summary's frequency distribution to the original source text frequency distribution (Reeve, Han, Nagori et al., 2006). Each frequency distribution (candidate summary and original source text) is modeled as a vector of unit items. Similarity functions are then applied to the two vectors. The five similarity functions used are: (1) Cosine similarity (Baeza-Yates & Ribeiro-Neto, 1999); (2) Dice's coefficient (Dice, 1945); (3) Euclidean distance; (4) vector subtraction (Subhash, 1996); and (5) vector model comparison considering only unit item frequency (Lee, Chuang, & Seamons, 1997). We found that Dice's coefficient, which looks at the number of common terms between the two vectors, performed the best (Reeve, Han, Nagori et al., 2006).

### 3.3. Hybrid method (ChainFreq)

The BioChain and FreqDist methods use different approaches to identify relevant sentences for building an extractive summary. A problem not addressed in the current BioChainSumm summarizer (Reeve, Han,

```
Initialization:
  INITIALIZE source-sentences to source-text sentences;
  INITIALIZE important-sentences to NULL;

Summary Generation:
  important-sentences = BioChain(source-sentences);

  important-sentences = FreqDist(source-text, important-sentences,
    summary size);

RETURN  important-sentences as final summary;
```

Fig. 4. Hybrid summarization method *ChainFreqSumm* using the *BioChain* method to identify sentences, and the *FreqDist* method to remove redundancy.

199  Nagori et al., 2006) is reducing information redundancy. Sentences containing the strongest concepts in the
200  text are extracted without a complimentary method for reducing redundancy from sentences already selected.
201  To overcome this limitation, we propose combining the BioChain and FreqDist methods to form a hybrid
202  method, called ChainFreq. The hybrid ChainFreq method first uses the BioChain method to identify candi-
203  date sentences containing strong concepts. The candidate sentences ($Sc$) and their corresponding concepts ($Cc$)
204  are then passed to the FreqDist method, which produces a set of summary sentences from the candidate sen-
205  tences. That is, only $Sc$ are used as a pool of sentences in the FreqDist method in Section 3.2. A summary
206  frequency distribution model is then created from the $Cc$, and their frequency counts are initialized to zero.
207  The FreqDist method then selects sentences containing concepts in the same distribution as the original source
208  text with respect to $Cc$ which reduces redundancy to the same proportion it exists in the source text.
209      Fig. 4 shows how the two summarization methods, BioChain and FreqDist, are combined to form the new
210  hybrid summarizer, ChainFreqSumm. First, all source sentences with their corresponding concept annota-
211  tions are collected and passed to the BioChain method. The BioChain method takes advantage of domain-spe-
212  cific knowledge, specifically UMLS semantic types, to find sentences which are important in the domain. There
213  is no limit on the number of sentences generated by the BioChain method. The subset of source-text sentences
214  identified by the BioChain method are then passed to the FreqDist method. The FreqDist method then finds a
215  further subset of sentences whose concept distribution best aligns with the concept distribution of the source
216  text. A user-defined summary size limits the number of sentences output at this stage. Both the BioChain
217  method and the FreqDist method work together to (a) find the important sentences according to the domain
218  (using the BioChain method), and (b) reduce redundancy by further reducing the number of important sen-
219  tences based on how well their concept distribution aligns with the source text's concept distribution (using the
220  FreqDist method), which has the effect of reducing redundancy.

221  **4. Evaluation**

222      The purpose of the evaluation is to evaluate the usefulness of concept frequency as a singular feature for
223  identifying salient sentences for extractive text summarization. The evaluation was done by first asking three
224  domain experts to manually generate extractive summaries from 24 biomedical texts (see Section 4.1). A series
225  of automated summarizers (see Section 4.5) then generated summaries of the biomedical texts. The output of
226  each summarizer is automatically compared using an automated tool called ROUGE (Lin, 2005) (see Section
227  5.3). The results are given in Section 5. The rest of this section details the evaluation implementation.

228  *4.1. Corpus*

229      A corpus of 24 biomedical texts was generated from a citation database of oncology clinical trial papers.
230  The database contains approximately 1200 papers physicians feel are important to the field (Brooks & Suli-
231  manoff, 2002). Of the 1200 papers cited, 24 were randomly selected. The PDF versions of these papers were
232  then obtained and converted to plain-text format. The papers were manually processed to remove graphics,
233  tables, figures, captions, citation references, and the bibliography section. The resulting text was further split

234 into an abstract text and a full-text source text (without the abstract). The number of papers chosen (24) was
235 based on the minimum requirements of the ROUGE summary evaluation tool (Lin, 2004) as well as the
236 resources available to complete the manual processing of each paper.

### 4.2. Rouge

238    The ROUGE (Recall-Oriented Understudy for Gisting Evaluation ) tool (version 1.5.5) (Lin & Hovy, 2003)
239 developed by the Information Science Institute at the University of Southern California was used. ROUGE
240 is an automated tool which compares a generated summary from an automated system with one or more
241 ideal summaries. The ideal summaries are called models. ROUGE uses N-grams to determine the overlap
242 between a summary and the models. ROUGE was used in the 2004 and 2005 Document Understanding
243 Conferences (DUC) (National Institute of Standards and Technology (NIST), 2005) as the evaluation tool.
244 We used parameters from the DUC 2005 conference. ROUGE-2 and ROUGE-SU4 recall scores are used
245 to measure each summarizer. ROUGE-2 evaluates bigram co-occurrence while ROUGE-SU4 evaluates ''skip
246 bigrams,'' which are pairs of words (in sentence order) having intervening word gaps no larger than four
247 words.

### 4.3. Model summaries

249    To compare summaries generated automatically from systems, we used four models (i.e., four ideal sum-
250 maries) for each of the 24 papers. The models represent different versions of ideal summaries. The first model
251 is the abstract of the paper (author's summary). In addition, three models from three different domain experts
252 were generated. The domain experts are medical students in their final year. Each was given the task of per-
253 forming extractive text summarization by selecting 20% of the sentences within a paper which formed the best
254 summary for that paper. Selecting a summary size was problematic. The news summarization domain typi-
255 cally selects a size of less than five sentences, which represents about 20% of the size of a typical news story
256 (Goldstein, Kantrowitz, Mittal, & Carbonell, 1999). It has been generally thought that a summary should be
257 no shorter than 15% and no longer than 35% of the source text (Hovy, 2005).

### 4.4. Additional summarizers

259    In this evaluation, eight additional extractive summarizers were randomly selected based on the type of
260 summarization method and availability. There are roughly four categories of summarizers selected: baseline,
261 frequency-based, multiple feature, and redundancy-sensitive, and we select two summarizers in each category.
262 The two baseline summarizers are Baseline-Lead, which sequentially selects the first 20% of sentences in the
263 source text, and Baseline-Random, which randomly selects 20% of the sentences in the source text. The fre-
264 quency-based summarizers are AutoSummarize in Microsoft Word (Microsoft Coporation, 2002) and Open
265 Text Summarizer (OTS) (Rotem, 2003). AutoSummarize is a feature of the Microsoft Word (Microsoft Copo-
266 ration, 2002) word processing software, and although exact details of the algorithm are not documented,
267 online help for the product indicates sentences using frequently-used words are given a higher score than sen-
268 tences containing low frequency words. OTS is an open source project where stemming can also be performed
269 to eliminate word variations. The two summarizers using multiple features to identify sentences are SweSum
270 (Dalianis, 2000) and MEAD (Radev et al., 2004). SweSum is a multi-lingual summarizer for Swedish and Eng-
271 lish text using features such as sentence position and numerical data identification. MEAD is a single and
272 multi-document summarizer using features such as position of sentence within the text, overlap of each sen-
273 tence with the first sentence, sentence length, and a centroid method based on a cluster of related documents.
274 Finally, the two summarizers which reduce information redundancy are Lemur Maximal Marginal Relevance
275 (MMR) (The Lemur Project, 2006) and SumBasic (Nenkova & Vanderwende, 2005). Lemur MMR iteratively
276 selects sentences having a high query similarity to an automatically-generated query, and which are also max-
277 imally dissimilar to sentences already included in the summary. SumBasic uses a probability distribution of
278 terms in the text, and reduces term probability as sentences containing the terms are selected. SumBasic
279 was also adapted to use concepts as the input source, rather than terms.

## 5. Results

### 5.1. Summarization method performance

The ROUGE results for all summarizers is shown in Tables 1 and 2. The following sections describe the ROUGE performance of the various summarization methods.

#### 5.1.1. BioChainSumm summarizer

For the BioChainSumm summarizer, the most frequent strong chain concept (MostFrequentStrongChain-Concept) scoring method outperforms the use of all strong chain concepts (AllStrongChainConcepts) according to both ROUGE-2 and ROUGE-SU4. The use of AllStrongChainConcepts does not have the effect of finding the most important sentences, but instead a broad coverage of the text, since all concepts within a strong chain are used to score sentences. While at first the broad coverage may seem desirable for summarizing a complete text, the result is that this technique does not find the most important sentences, which is what is desired in a summary. By using the MostFrequentStrongChainConcept, the effect is to use the UMLS semantic type to find the general idea of a text, and then the main concept within the chain to further refine the idea. Both versions of BioChainSumm outperform the two baseline summarizers Baseline-Lead and Baseline-Random, as

Table 1
ROUGE-2 scores

| Summarizer | ROUGE-2 score |
| --- | --- |
| FreqDistSumm-Term-Dice | 0.12653 |
| ChainFreqSumm-AllStrongChainConcepts-Dice | 0.12216 |
| FreqDistSumm-Concept-Dice | 0.12070 |
| SumBasic-Term | 0.11673 |
| SumBasic-Concept | 0.10940 |
| Lemur-MMR | 0.10708 |
| ChainFreqSumm-MostFrequentStrongChainConcept-Dice | 0.10652 |
| BioChainSumm-MostFrequentStrongChainConcept | 0.10419 |
| BioChainSumm-AllStrongChainConcepts | 0.09708 |
| Mead | 0.09254 |
| Baseline-Random | 0.08001 |
| MSWord | 0.07977 |
| SweSum | 0.07513 |
| OTS | 0.07474 |
| Baseline-Lead | 0.07076 |

Table 2
ROUGE-SU4 scores

| Summarizer | ROUGE-SU4 score |
| --- | --- |
| ChainFreqSumm-AllStrongChainConcepts-Dice | 0.22303 |
| FreqDistSumm-Term-Dice | 0.22176 |
| FreqDistSumm-Concept-Dice | 0.21997 |
| SumBasic-Term | 0.21112 |
| ChainFreqSumm-MostFrequentStrongChainConcept-Dice | 0.20158 |
| SumBasic-Concept | 0.20034 |
| Lemur-MMR | 0.19874 |
| BioChainSumm-MostFrequentStrongChainConcept | 0.19173 |
| BioChainSumm-AllStrongChainConcepts | 0.18557 |
| Mead | 0.17629 |
| Baseline-Random | 0.16396 |
| MSWord | 0.15171 |
| SweSum | 0.15115 |
| OTS | 0.14919 |
| Baseline-Lead | 0.13953 |

294  well as the generic summarizers AutoSummarize, OTS, SweSum, and MEAD. This indicates the use of domain-
295  specific concepts for selecting sentence is an improvement over the use of terms which lack domain-specific
296  knowledge. However, neither version of BioChainSumm outperforms SumBasic, Lemur (MMR), FreqDist-
297  Summ, nor the hybrid ChainFreqSumm. All of these summarizers have a redundancy-removal component
298  as part of their design. This shows that while BioChainSumm is effective at identifying important sentences
299  within a domain, it falls short of removing redundancy, which is an important part of text summarization.

### 5.1.2. FreqDistSumm summarizer

301  The term and concept versions of FreqDistSumm outperform all other summarizers except for the hybrid
302  ChainFreqSumm (BioChain method plus the FreqDist method) (see Section 5.1.3). In both ROUGE-2 and
303  ROUGE-SU4, the term version of FreqDistSumm outperforms the concept version. Interestingly, this is also
304  true for the SumBasic summarizer. We believe this is due to the fact that concepts map multiple instances of an
305  expression to a single concept, and in single document summarization there is not enough variation in lan-
306  guage to allow concepts to outperform terms. For example, in two articles discussing lung cancer, one article
307  may use 'lung cancer' repeatedly, whereas the second article uses 'pulmonary carcinoma.' Using concepts,
308  these two instances will be merged. When using terms, the different instances will remain unique. The FreqDist
309  method keeps redundancy in check by only allowing it to occur in the same degree as the source text. Reiter-
310  ation is a technique often used by authors to emphasize important points. The FreqDistSumm summarizer
311  attempts to mirror the source-text in a reduced form, and so the reiteration will also be expressed in the gen-
312  erated summary. This is a different approach than BioChainSumm, which finds all important sentences using
313  domain-defined criteria, but does not then reduce any redundancy in the important sentences. In contrast,
314  FreqDistSumm will find a subset of source-text sentences using domain-defined criteria, but then eliminates
315  some of the sentences from the subset if they emphasize main points more than the source text does.

### 5.1.3. ChainFreqSumm summarizer

317  As can be seen from the Baseline-Random summarizer, randomly picking sentences performs well. This is
318  an indication that biomedical texts contain a large amount of redundancy. As can also be seen from the Bio-
319  ChainSumm evaluation (see Section 5.1.1), redundancy can decrease summarizer performance. The hybrid
320  ChainFreqSumm summarizer (BioChain method plus the FreqDist method) is an attempt to find a subset
321  of the most important sentences using domain-specific criteria, and then remove redundancy from the subset.
322  The ChainFreqSumm summarizer performs best when all concepts in the strong chains are used, which is the
323  opposite of what occurs when the BioChain method is used alone. This is most likely because using all strong
324  concepts results in a larger pool of sentences for the FreqDist method to select from. Using the ROUGE-SU4
325  metric, the hybrid ChainFreqSumm summarizer is the best performer, but is slightly outperformed by the
326  FreqDistSumm term method when the ROUGE-2 metric is used. The result in combining the two approaches
327  is that the use of concept approaches for finding salient sentences is improved over the individual methods of
328  FreqDist and BioChain. We believe that a summarizer which (a) first identifies a subset of important sentences
329  based on domain-specific criteria, and (b) then prunes the subset by removing redundancy leads to an effective
330  domain-specific summarizer.

### 5.2. Physician evaluation

332  We provided a practicing oncologist our corpus of evaluation papers and asked him to select several papers
333  and give his observations on the abstract and the full-text. He chose three papers based on their chaining per-
334  formance in the original BioChainSumm evaluation (Reeve, Han, Nagori et al., 2006). The intent is to sub-
335  jectively evaluate the usefulness of abstracts of a full-text in the practice of oncology treatment. The
336  following observations were made about each paper:
337  Paper #1 (Riethmuller et al., 1998):

338  – The abstract and the full-text had different information regarding the number of patients.
339  – Information about the study design was not complete in the abstract, causing a misleading conclusion.
340  Paper #2 (Perry et al., 1987):

341 – The abstract gives conclusions without providing context of the results, such as follow-up time and
342 response criteria.
343    Paper #3 (Thomas et al., 1998):
344 – The abstract makes a statement that is not made in the source text.
345

346    The physician concludes that the abstracts are a good starting point, but that they may miss critical infor-
347 mation necessary for evaluating the results or overstate the conclusion. The physician also believes some bias
348 is inherent in author-generated abstracts, and that additional ways of producing summaries are needed to
349 address different needs of medical practitioners.
350    We followed up the physician's evaluation to see if our FreqDistSumm-generated summaries addressed the
351 shortcomings of the abstracts. For Paper #1, complete study design information was included. In Paper #2,
352 the summary included additional context information, such as patient eligibility and stratification and ran-
353 domization. In Paper #3, the full-text did not contain the information in the abstract and so could not possibly
354 be extracted. Based on this short evaluation, we note that automatically generated summaries addressed
355 abstract incompleteness in two out of three cases.

356 **6. Conclusion**

357    We presented three novel semantic-based methods for extractive text summarization. The first method (Bio-
358 Chain) chains together semantically-related concepts, and then extracts sentences having concepts in the stron-
359 gest of the chains. The second method (FreqDist) uses a frequency-distribution approach, where a summary is
360 gradually constructed by adding new source sentences so that the summary and source text have similar con-
361 cept frequency distributions. For single document summarization, we show that these two concept-based
362 approaches are competitive with existing term-based approaches. In addition, we combine the two approaches
363 (BioChain and FreqDist) to improve performance above all existing summarizers. The use of concepts can be
364 more useful than terms for generating personalized summaries and multi-document summarization. An envi-
365 sioned system allows a user to select domain-specific concepts important to the user, and then have the sum-
366 marizer generate a summary where those concepts are more highly weighted than the concepts appearing in
367 the source text.

368 **References**

369 Baeza-Yates, R., & Ribeiro-Neto, B. (1999). *Modern information retrieval* (1st ed.). Harlow, England: Addison-Wesley.
370 Barzilay, R., & Elhadad, M. (1997). Using lexical chains for text summarization. In *Proceedings of the intelligent scalable text*
371    *summarization workshop (ISTS'97)* (pp. 10–18). ACL, Madrid, Spain.
372 Brooks, A. D., & Sulimanoff, I. (2002). Evidence-based oncology project. *Surgical Oncology Clinics of North America*, 3–10.
373 Carbonell, J., & Goldstein, J. (1998). The use of MMR, diversity-based reranking for reordering documents and producing summaries. In
374    *SIGIR '98: Proceedings of the 21st annual international ACM SIGIR conference on research and development in information retrieval* (pp.
375    335–336). Melbourne, Australia. Available from http://doi.acm.org/10.1145/290941.291025.
376 Cohen, A. M., & Hersh, W. R. (2005). A survey of current work in biomedical text mining. *Briefings in Bioinformatics, 6*(1), 57–71.
377 D'Avanzo, E., Magnini, B., & Vallin, A. (2004). Keyphrase extraction for summarization purposes: The LAKE system at DUC-2004. In
378    *Proceedings of the 2004 document understanding conference*, Boston, USA.
379 Dalianis, H. (2000). *SweSum – A text summarizer for Swedish No. TRITA-NA-P001(5)*. Stockholm, Sweden: NADA, KTH.
380 Dice, L. R. (1945). Measures of the amount of ecologic association between species. *Ecology, 26*, 297–302.
381 Doran, W. P., Stokes, N. S., Dunnion, J., & Carthy, J. (2004). Assessing the impact of lexical chain scoring methods and sentence
382    extraction schemes on summarization. In *Proceedings of the 5th international conference on intelligent text processing and computational*
383    *linguistics CICLing-2004* (Vol. 2945, pp. 627–635). Seoul, South Korea.
384 Edmundson, H. P. (1999). New methods in automatic extracting. In I. Mani & M. T. Maybury (Eds.) (pp. 23–42). Cambridge, MA: MIT
385    Press.
386 Fellbaum, C. (1998). *WORDNET: An electronic lexical database*. Cambridge, MA: The MIT Press.
387 Galley, M., & McKeown, K. (2003). Improving word sense disambiguation in lexical chaining. In *Proceedings of the 18th international joint*
388    *conference on artificial intelligence* (pp. 1486–1488). Acapulco, Mexico.
389 Goldstein, J., Kantrowitz, M., Mittal, V., & Carbonell, J. (1999). Summarizing text documents: sentence selection and evaluation metrics.
390    In *SIGIR '99: Proceedings of the 22nd annual international ACM SIGIR conference on research and development in information retrieval*
391    (pp. 121–128). Berkeley, California, United States. Availbale from http://doi.acm.org/10.1145/312624.312665.

392  Hovy, E. H. (2005). Automated text summarization. In R. Mitkov (Ed.), *The oxford handbook of computational linguistics* (pp. 583–598).
393      Oxford: Oxford University Press.
394  Hovy, E., & Lin, C. (1999). Automated text summarization in SUMMARIST. In I. Mani & M. T. Maybury (Eds.), *Advances in automatic*
395      *text summarization* (pp. 81–94). Cambridge, MA: MIT Press.
396  Jaques, D. P. (Ed.). (2002). *Surgical oncology clinics of North America: Prospective randomized clinical trials in oncology* (1st ed.).
397      Philadelphia, PA, USA: W.B. Saunders Company.
398  Lee, D. L., Chuang, H., & Seamons, K. (1997). Document ranking and the vector-space model. *Software, IEEE, 14*(2), 67–75.
399  Lin, C. (2004). Looking for a few good metrics: Automatic summarization evaluation – How many samples are enough? In *Proceedings of*
400      *the NTCIR workshop 4*, Tokyo, Japan.
401  Lin, C., & Hovy, E. H. (2003). Automatic evaluation of summaries using N-gram co-occurrence statistics. In *Proceedings of 2003 language*
402      *technology conference (HLT-NAACL 2003)* (Vol. 1(1), pp. 71–78). Edmonton, Canada.
403  Lin, C. (2005). Recall-oriented understudy for gisting evaluation (ROUGE). Retrieved August 20, 2005, from http://www.isi.edu/~cyl/
404      ROUGE/.
405  Luhn, H. P. (1958). The automatic creation of literature abstracts. *IBM Journal of Research and Development, 2*(2), 159–165.
406  Microsoft Coporation. (2002). *Microsoft word 2002*. Redmond, Washington, USA.
407  Morris, J., & Hirst, G. (1991). Lexical cohesion computed by thesaural relations as an indicator of the structure of text. *Computational*
408      *Linguistics, 17*(1), 21–43.
409  National Institute of Standards and Technology (NIST). (2005). Document understanding conferences. Retrieved August 20, 2005, from
410      http://www-nlpir.nist.gov/projects/duc/.
411  Nenkova, A., & Vanderwende, L. (2005). *The impact of frequency on summarization No. MSR-TR-2005-101*. Redmond, Washington:
412      Microsoft Research.
413  Perry, M. C., Kardinal, C. G., Korzun, A. H., Ginsberg, S. J., Raich, P. C., Holland, J. F., et al. (1987). Chemohormonal therapy in
414      advanced carcinoma of the breast: Cancer and Leukemia Group B protocol 8081. *Journal of Clinical Oncology: Official Journal of the*
415      *American Society of Clinical Oncology, 5*(10), 1534–1545.
416  Pollock, J. J., & Zamora, A. (1975). Automatic abstracting research at chemical abstracts service. *Journal of Chemical Information and*
417      *Computer Sciences, 15*(4), 226–232.
418  Radev, D., Allison, T., Blair-Goldensohn, S., Blitzer, J., Celebi, A., & Drabek, E. et al. (2004). MEAD – A platform for multidocument
419      multilingual text summarization. In *Proceedings of language resources and evaluation 2004 (LREC 2004)*. Lisbon, Portugal.
420  Rath, G. J., Resnick, A., & Savage, R. (1961). The formation of abstracts by the selection of sentences [Electronic version]. *American*
421      *Documentation, 2*(12), 139–208.
422  Reeve, L., Han, H., & Brooks, A. D. (2006). BioChain: Using lexical chaining methods for biomedical text summarization. In *Proceedings*
423      *of the 21st annual ACM symposium on applied computing, bioinformatics track* (pp. 180–184). Dijon, France.
424  Reeve, L., Han, H., Nagori, S. V., Yang, J., Schwimmer, T., & Brooks, A. D. (2006). Concept frequency distribution in biomedical text
425      summarization. In *Proceedings of the ACM 15th conference on information and knowledge management (CIKM'06)*. Arlington, VA,
426      USA.
427  Riethmuller, G., Holz, E., Schlimok, G., Schmiegel, W., Raab, R., Hoffken, K., et al. (1998). Monoclonal antibody therapy for resected
428      Dukes' C colorectal cancer: Seven-year outcome of a multicenter randomized trial. *Journal of Clinical Oncology: Official Journal of the*
429      *American Society of Clinical Oncology, 16*(5), 1788–1794.
430  Rotem, N. (2003). Open text summarizer (OTS). Retrieved July 3, 2006, from http://libots.sourceforge.net.
431  Silber, G. H., & McCoy, K. F. (2002). Efficiently computed lexical chains as an intermediate representation for automatic text
432      summarization. *Computational Linguistics, 28*(4), 487–496.
433  Sparck Jones, K. (1999). Automatic summarizing: Factors and directions. In I. Mani & M. T. Maybury (Eds.), *Advances in automatic text*
434      *summarization* (pp. 2–12). Cambridge, MA: MIT Press.
435  Subhash, S. (1996). *Applied multivariate techniques* (1st ed.). USA: John Wiley and Sons.
436  The Lemur Project. (2006). Lemur Language Modeling Toolkit. Retrieved July 3, 2006, from http://www.lemurproject.org/.
437  Thomas, G., Dembo, A., Ackerman, I., Franssen, E., Balogh, J., Fyles, A., et al. (1998). A randomized trial of standard versus partially
438      hyperfractionated radiation with or without concurrent 5-fluorouracil in locally advanced cervical cancer. *Gynecologic oncology, 69*(2),
439      137–145.
440  United States National Library of Medicine. (2006a). PubMed. http://www.ncbi.nlm.nih.gov/entrez/query.fcgi.
441  United States National Library of Medicine. (2006b). UMLS Metathesaurus fact sheet. http://www.nlm.nih.gov/pubs/factsheets/
442      umlsmeta.html.
443  United States National Library of Medicine. (2005a). ClinicalTrials.gov. http://www.clinicaltrials.gov/.
444  United States National Library of Medicine. (2005b). MetaMap transfer. http://mmtx.nlm.nih.gov/.
445  United States National Library of Medicine. (2005c). Unified medical language system (UMLS). http://www.nlm.nih.gov/research/umls/.
446  United States National Library of Medicine. (2004). UMLS semantic network fact sheet. http://www.nlm.nih.gov/pubs/factsheets/
447      umlssemn.html.
448