

[College of Information Science and Technology](#)



Drexel E-Repository and Archive (iDEA)
<http://idea.library.drexel.edu/>

Drexel University Libraries
www.library.drexel.edu

The following item is made available as a courtesy to scholars by the author(s) and Drexel University Library and may contain materials and content, including computer code and tags, artwork, text, graphics, images, and illustrations (Material) which may be protected by copyright law. Unless otherwise noted, the Material is made available for non profit and educational purposes, such as research, teaching and private study. For these limited purposes, you may reproduce (print, download or make copies) the Material without prior permission. All copies must include any copyright notice originally included with the Material. **You must seek permission from the authors or copyright owners for all uses that are not allowed by fair use and other provisions of the U.S. Copyright Law.** The responsibility for making an independent legal assessment and securing any necessary permission rests with persons desiring to reproduce or use the Material.

Please direct questions to archives@drexel.edu

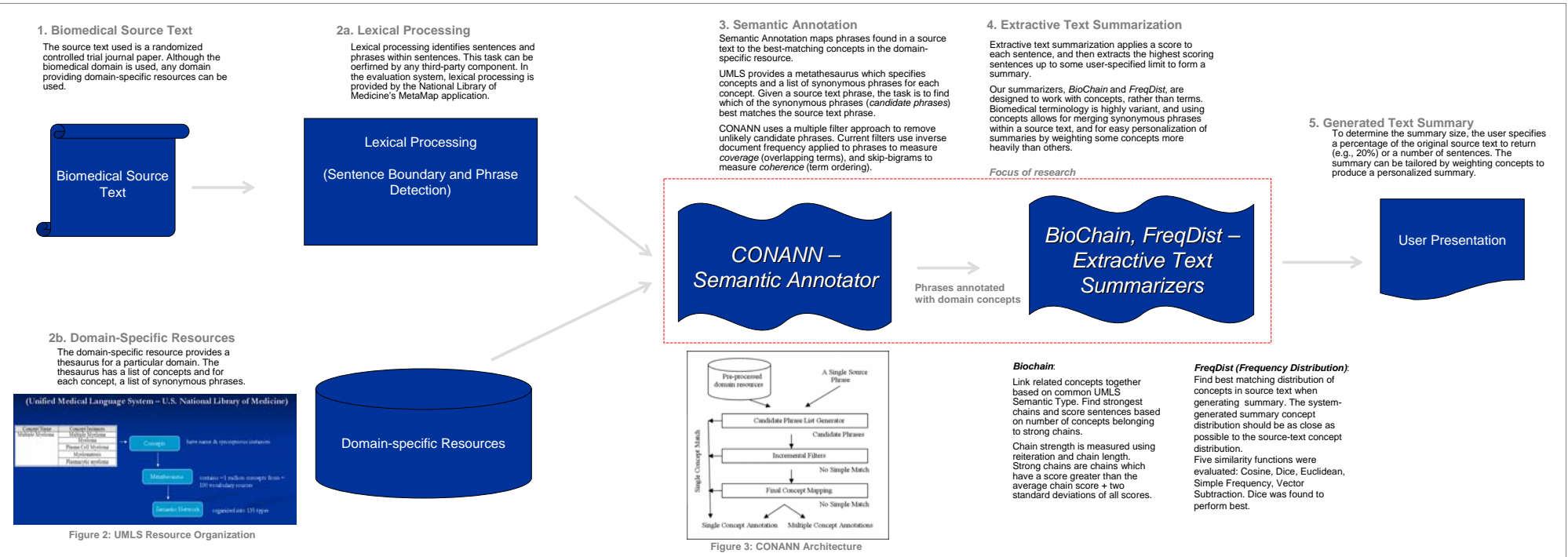
Biomedical Text Annotation and Summarization

Lawrence H. Reeve¹, Hoyil Han¹, Ari D. Brooks²

¹ Drexel University, College of Information Science and Technology

² Drexel University, College of Medicine

Figure 1: Text Summarization System and Semantic Annotation system



Motivation, Hypothesis and Method

Motivation: Generate short summaries of biomedical texts (randomized controlled trials in oncology) to allow physicians and researchers to assimilate more information in less time.

Approach: Identify and extract sentences from the source text to form a summary. The problem is how to identify important sentences within the full source text to extract while simultaneously reducing information redundancy.

Hypothesis: Domain-specific concepts can be used to identify important areas with a text which can then be extracted to form a summary.

Method: Annotate biomedical source texts with biomedical concepts and then use two novel algorithms which utilize concepts rather than terms to identify sentences which express the main themes of the text.

Publications

Reeve, L., Han, H., & Brooks, A. D. (2006). Biomedical Text Summarization Using Concept Chains. *International Journal of Data Mining and Bioinformatics*, 1(4), 389-407.

Reeve, L. H., Han, H., & Brooks, A. D. (2007). The Use of Domain-Specific Concepts in Biomedical Text Summarization. To appear in the *Journal of Information Processing and Management, Special Issue on Summarization*. Elsevier.

Reeve, L. H. & Han, H. (2007). CONANN: An Online Biomedical Semantic Annotator. To appear in the *Proceedings of Data Integration in the Life Sciences*. June 27-29, 2007.

Summarization Evaluation

Basic Idea

- Generate model summaries from domain experts
- Generate system summaries
- Use ROUGE to compare system and model summaries

Semantic Annotation Evaluation

Intrinsic

- Use NLM's MetaMap application to generate baseline of singly-mapped noun phrases
- Pass this set of phrases to CONANN
- Compare precision and time-to-annotate

Extrinsic

- Compare MetaMap vs. CONANN in text summarization performance using FreqDist and SumBasic summarizers

Observations

- CONANN Annotator:
 - Performs 20x faster than state-of-the-art system
 - Small loss of precision, acceptable in a user task
- BioChain, FreqDist Summarizers:
 - Concepts perform closely with terms
 - Concepts can be used to personalize summaries

Abstract A two-part text summarization system is described which addresses the scalability of utilizing the large volume of text-based information in the biomedical field. The contributions of the system are a) it utilizes biomedical concepts rather than terms to find the main points of a text, b) it uses two new concept-based algorithms to find important areas of a text for extracting sentences to form a summary, and c) it is supported by a new semantic annotation sub-system, which identifies biomedical concepts found in biomedical text documents. The semantic annotation sub-system uses a novel multiple-filter system architecture for online matching of concepts defined by a biomedical metathesaurus. The goal of semantic annotation is to show online text-to-concept mapping can be performed without a significant loss of precision as compared to current offline systems. An evaluation shows the text summarization algorithms using concepts outperform existing summarization systems, and the semantic annotation system performs twenty times faster than a state-of-the-art system with no significant loss of precision.

Semantic Annotation Results

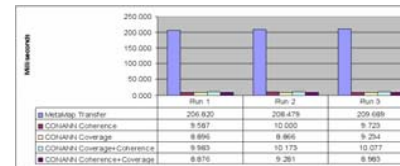


Figure 4: Average Phrase Annotation Time



Figure 5: Annotator Precision

Summarization Method	ROUGE-3 Score	ROUGE-SU4 Score
FreqDist using MetaMap	0.1207	0.2200
FreqDist using CONANN	0.1192	0.2161
SumBasic using CONANN	0.1178	0.2098
SumBasic using MetaMap	0.1094	0.2003

Figure 6: ROUGE Scores (extrinsic evaluation)

Text Summarization Results

Summarizer	ROUGE-2	Summarizer	ROUGE-SU4
FreqDist Term Dice	0.1371	FreqDist Term Dice	0.2005
SumBasic Term	0.1271	SumBasic Term	0.2169
FreqDist ConceptBMMTX Dice	0.1260	FreqDist ConceptBMMTX Dice	0.2164
FreqDist ConceptBMMTX Dice	0.1153	FreqDist ConceptBMMTX Dice	0.2107
SumBasic ConceptBMMTX	0.1157	SumBasic ConceptBMMTX	0.2084
SumBasic ConceptBMMTX	0.1094	SumBasic ConceptBMMTX	0.2077
SumBasic ConceptBMMTX	0.1072	SumBasic ConceptBMMTX	0.1910
SumBasic ConceptBMMTX	0.1061	SumBasic ConceptBMMTX	0.1872
SumBasic ConceptBMMTX	0.1041	SumBasic ConceptBMMTX	0.1810
SumBasic ConceptBMMTX	0.1026	SumBasic ConceptBMMTX	0.1774
SumBasic ConceptBMMTX	0.1010	SumBasic ConceptBMMTX	0.1700
SumBasic ConceptBMMTX	0.0984	SumBasic ConceptBMMTX	0.1686
SumBasic ConceptBMMTX	0.0974	SumBasic ConceptBMMTX	0.1645
SumBasic ConceptBMMTX	0.0934	SumBasic ConceptBMMTX	0.1464
SumBasic ConceptBMMTX	0.0711	SumBasic ConceptBMMTX	0.1376
SumBasic ConceptBMMTX	0.0291	SumBasic ConceptBMMTX	0.0830
SumBasic ConceptBMMTX	0.0143	SumBasic ConceptBMMTX	0.0271

Figure 7: ROUGE-2 and ROUGE-SU4 scores for text summarization

Summarizer	Group	Count	Sum	Average	Variance
SumBasic	Group 1	9	0.81348	0.0903111	0.000264
SumBasic	Group 2	3	0.73346	0.2444867	0.28101

Figure 8: ANOVA: FreqDist vs. Other summarizers