

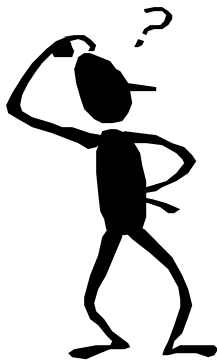
Ou, S., Khoo, C. & Goh, D. (2006). Multi-document summarization for digital libraries.
Presented at the Asia-Pacific Conference on Library & Information Education & Practice 2006 (A-LIEP 2006), Singapore, 3-6 April 2006, Nanyang Technological University.

Multi-document Summarization for Digital Libraries

Ou Shiyan, Khoo Christopher, Goh Dion
Division of Information Studies
School of Communication & Information
Nanyang Technological University

1

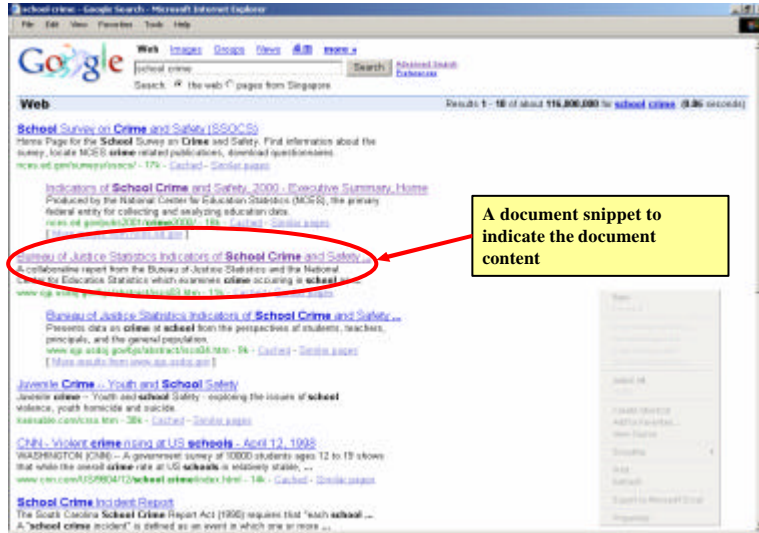
Background



- **How to solve the problem of information overload?**
 - Display a list of titles and short abstracts of retrieved documents (e.g. [Google](#));
 - Group the retrieved documents into folders or categories (e.g. [Northern Light search engine](#));
 - **Construct a multi-document summary of the documents retrieved.**

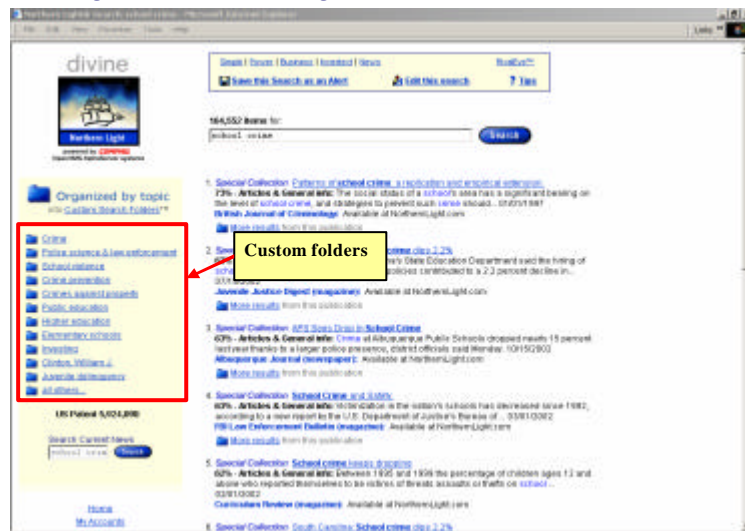
2

A list of document “snippets” returned by Google



3

Search results clustered into custom folders by the North Light search engine



4

Research Objectives

- Develop a method for automatic construction of multi-document summaries of sets of related research abstracts;
- Focus on **research concepts (or variables)** and **relationships**;
- Carry out a user evaluation to compare the variable-based summary against two sentence-based summaries.

5

Automatic Summarization

- Process of creating a shorter representation of an original information source;
- Different types of summaries:
 - Extracts vs. Abstracts
 - Indicative vs. Informative summaries
 - Generic vs. User-focused summaries
 - Single-document vs. Multi-document summaries

6

Automatic Multi-document Summarization

- More useful in digital libraries;
- More challenges in the issues of compression, redundancy, cohesion, coherence, and so on;
- New approaches for multi-document summarization
 - Reduce redundancy;
 - Ensure cohesion and coherence;
 - Identify similarities and differences across documents.

7

Summarization Approaches

- Extract important sentences
 - **Strengths:** domain independent; don't need understand the meaning and structure of the text;
 - **Weaknesses:** the resulting extracts are redundant, not coherent and fluent, and hard to read;
- Identify similarities and differences among a set of documents
 - **Similarities** indicate salient information in a document set;
 - **Differences** indicate unique information in individual documents;

8

Three Types of Multi-document Summaries

- **MEAD summary:** a sentence-based summary, generated by a state-of-the-art summarization system that extracts the highly ranked sentences across documents;
- **Research objective summary:** a sentence-based summary, generated by extracting research objective sentences from each document;
- **Variable-based summary:** a summary focusing on research concepts often operationalized variables, generated by extracting research concepts and relationships from each document and integrating them across documents.

9

MEAD Summary

- Use a cross-document sentence extraction approach;
- Generated by a state-of-the-art multi-document summarization system, MEAD;
- MEAD, domain public, built by the University of Michigan;
- Cluster similar documents and extract the most important sentences from each cluster using various features;

10

A MEAD summary on the topic of “school crime”



11

Research Objective Summary

- Use sentence extraction approach;
- Extract research objective sentences from each dissertation abstract;
- Identify research objective sentences using a decision tree classifier and cue phrases found at the beginning of some sentences, such as “*The purpose of the study was to explore...*”, “*This study was to investigate ...*”.

12

A research objective summary on the topic of “school crime”

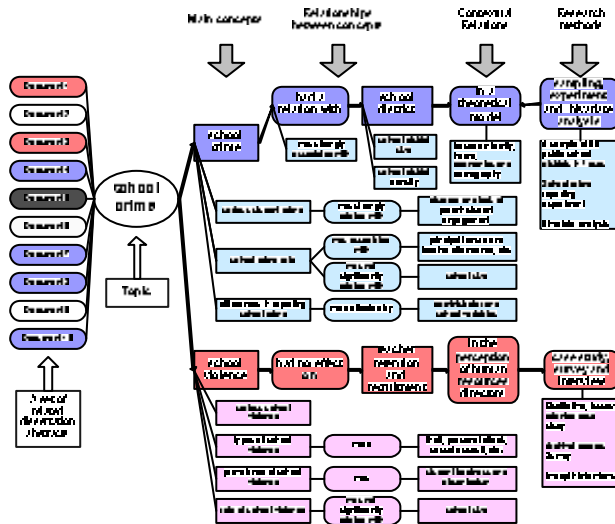
The screenshot shows a web browser window with the title "OBJECTIVES summary on the topic of 'school crime'" and a subtitle "(all the research objectives sentences from 18 dissertation abstracts)". The main content is a list of 18 numbered bullet points, each representing a research objective from a different dissertation abstract. The objectives cover various aspects of school crime, including teacher retention, legal definitions of violence, prevalence in specific districts, school administrator differences, school size and dropout rates, restructuring impacts, theoretical models, and factors related to student-initiated crimes.

13

Variable-based Summary

- Focus on research concepts and relationships;
- Much of sociology research aims to explore research concepts and relationships (Macionis, 2000).
- The similarities and differences across dissertation abstracts are mainly reflected through **research concepts** and **relationships**.
- A variable-based framework is proposed for integrating and organizing information extracted from different dissertation abstracts and thus summarizing a set of dissertation abstracts.

Variable-based Framework



15

Relationships associated with “school crime” in 5 of 10 dissertation abstracts

- **School crime** was strongly associated with school district size and density;
- School size was not significantly related with school crime rate;
- **Serious student crime** was strongly related with disorder and lack of parent/student engagement;
- **School crime rate** was associated with principal tenure and teacher attendance, abolition of corporal punishment etc;
- **Differences in reporting school crime** was affected by school and administrator variables.

16

Integrated concepts and relationships associated with “school crime”

- **School crime**, including **serious student crime**, **school crime rate**, and so on;

Its different aspects were investigated, including **differences in reporting school crime**, and so on.

The following relationships were investigated:

- It was related with **school district size and density**, **principal tenure** and **teacher attendance**, **abolition of corporal punishment**, **disorder** and **lack of parent/student engagement** etc;
- It was not related with **school size**;
- It was affected by **school** and **administrator variables**;

17

A variable-based summary on the topic of “school crime”

SYSTEM 2 summary on the topic of "school crime"

* Number in the brackets indicates the number of documents.
* Concepts highlighted in red are the more common concepts in sociology dissertation abstracts.

In these 10 dissertation abstracts, the following **related relations** were found:
 perception(2), attitude(1), framework(1), hypothesis(1), model(1), perspective(1), theory(1), view(1)

In these 10 dissertation abstracts, the following **research methods** were found:
 survey(1), case study(2), regression analysis(2), sampling(2), archival research(1), literature analysis(1), correlational research(1), empirical research(2), experiment(1), interview(1), multilevel analysis(1), multivariate analysis(1), etc.

These 10 dissertation abstracts were mainly about:

1. Education 2. Social and Human Sciences 3. Ethics and Law
 4. Economics 5. General Concepts

1. Education:
school(10) including public school(1), high school(1), student and school(2), daily school(1), individual school(1), large school(1), marching school(1), middle school(1), restructuring school(1), selected school(1), senior school(1), suburban school(1), school and community(1), school in the Canton(1), and more ...

Different aspects were investigated, including school crime(1), school district(1), school violence(2), school campus(2), school safety(2), school size(2), area of school(1), area of the school(1), aspect of school(1), behavior at school(1), commitment to school(1), disorder within the school(1), officer in school(1), reporting of school(1), and more ...

The following relationships were investigated:

- o There may be an effect on school delinquency .
- o There was no effect on selection and recruitment .
- o It was affected by number of areas for possible intervention by policy makers, administrator and school variables .
- o It may be affected by school architecture, faculty affairs .
- o There was a relation with number of students, school district, factors, principal tenure and teacher attendance, abolition of corporal punishment .

18

The screenshot displays a search results page for 'school crime'. On the left, a sidebar lists various categories: '2. Social and human sciences', '3. Politics and law', '4. Economics', and '5. General concepts'. Under '2. Social and human sciences', there are sub-categories: 'crime(2)', 'violence(1)', and 'assault(1)'. A red arrow points from the 'school crime(1)' link in the 'crime(2)' sub-category to the first document in the main list. The main list contains four documents, each with a title and a brief description of its content.

1. doc_id=4 [Difference in school administrator's reporting of school crime to law enforcement](#)
 school crime, school administrator reporting of school crime to law enforcement
 The research methods were positivistic, quantitative research.
 The research methods were five-category questionnaires, random sample.

2. doc_id=5 [The relationship of high school size to student extracurricular activity participation, student dropout rates, and crime and violence in the school](#)
 rate of school crime and violence and student extracurricular activity participation, school dropout rate and rate of school crime and violence, school dropout rate and school crime and rate of school crime and violence
 The research methods were not found.
 The research methods were multiple regression analysis.

3. doc_id=7 [School crime in Texas: an initial examination of a public problem](#)
 school crime, school crime in Texas, matching school crime data with various characteristics of the district and reporting factors, theoretical model of school crime, theoretical model to school crime
 The research methods were descriptive model of school crime, descriptive model to school crime.
 The research methods were logistic analysis of the district, reported data, regression, single factor regression equation.

4. doc_id=8 [Factors related to student-submitted serious crime on high school campus](#)
 student-submitted serious crime on high school campus, serious school crime, strong relationship to serious crime

19

Evaluation Design

- Compare the three types of summaries:
 - Variable-based summary
 - MEAD summary
 - Research objective summary
- 20 researchers in sociology domain as human subjects;
- Three types of summaries for 20 research topics submitted by these 20 researchers;
- Evaluate the three types of summaries on two criteria:
 - Quality: readability and comprehensibility
 - Usefulness for the research-related purpose

20

Evaluation Results (1)

- 14 researchers (70%) indicated their preference for the **variable-based summaries**;
- 11 researchers (55%) indicated their preference for the **research objective summaries**;
- 5 researchers (25%) indicated the preference for the **MEAD summaries**.

21

Evaluation Results (2)

- Variable-based summary:
 - efficient to give an overview of a topic;
 - well-organized and concise;
 - useful for information scanning;
 - too brief to provide accurate information;
 - have the potential to confuse users;
- Research objective summary:
 - more straight indication of the main points of the dissertation;
 - vague and indistinct;
- MEAD summary
 - provide a more complete overview of a topic;
 - hard to read.

22

Evaluation Results (3)

- The variable-based summary was found more useful in the following aspects:
 - Give an overview of a research area;
 - Indicate similarities among previous studies;
 - Indicate important concepts in the research area;
 - Indicate important research methods used in the research area;
- MEAD summary and research objective summary were found to be more useful in the following two aspects:
 - Help to identify the documents of interest easily;
 - Indicate important theories, views, or ideas in the research area.

23

Summary and Conclusion

- Introduce three types of summaries:
 - MEAD summary
 - Research objective summary
 - Variable-based summary
- Evaluate the three types of summaries
 - The majority of users (70%) preferred the variable-based summary to the sentence-based summaries
- Describe a new summarization method
 - Extract research concepts and relationships
 - Integrate similar concepts and same types of relationships
 - Use a concept-oriented presentation design

24

Thank you for your listening

Questions?

