



1st Workshop on Recommender Systems for Technology Enhanced Learning, RecSysTEL 2010

## Semantic Ranking of Lecture Slides based on Conceptual Relationship and Presentational Structure

Yuanyuan Wang<sup>a</sup>, Kazutoshi Sumiya<sup>b,\*</sup>

<sup>a</sup>Graduate School of Human Science and Environment, University of Hyogo, Japan

<sup>b</sup>School of Human Science and Environment, University of Hyogo, Japan

---

### Abstract

We describe a presentation content retrieval method involving the semantic ranking of target slides based on the relations between slides related to a user query. This method uses a keyword conceptual structure of the conceptual relationship implicitly existing between keywords extracted from the slide text of the presentation content and the presentational structure of indents in the slide text. At present, many presentation files are shared over the Web by many universities through their own public sites. Although these files are useful and valuable to many potential students, the fact that such files have to be retrieved for self-learning purposes means that there is still a lack of support for self-learners to find the desired slides on a priority basis, i.e. on the basis of importance and urgency of requirement of the desired file. Our noble semantic ranking method helps a user to easily learn through slides, focusing on either highly detailed slides or introductory slides in an order related to the user query. We also present a prototype system supported by our method for slide ranking and evaluate its effectiveness through experiments.

© 2010 Published by Elsevier B.V.

*Keywords:* Multimedia; E-learning; Presentation content retrieval; Slide ranking; Conceptual relationship

---

### 1. Introduction

Free online presentation contents often provide lecture slides. At present, a considerable amount of lecture material is shared on websites such as SlideShare<sup>1</sup> and MPMeister<sup>2</sup>. Thus, not only students who missed a lecture but also those interested in the topic being discussed in the lecture can review the lecture and study its content on their own at their convenience. When a user asks a query, he or she must know the query well in order to retrieve the required lecture slides on the basis of the matching keywords. If the keywords in a query tend to appear many times, there could be a possibility that many irrelevant slides will be retrieved. A simple retrieval method cannot retrieve relevant slides on the basis of a query; therefore, this method makes it difficult to obtain an appropriate retrieval result.

---

\* Corresponding author. Tel.: +81-792-92-9339; fax: +81-792-92-9339.

E-mail address: [sumiya@shse.u-hyogo.ac.jp](mailto:sumiya@shse.u-hyogo.ac.jp)

<sup>1</sup> <http://www.slideshare.net/>

<sup>2</sup> <http://www.ricoh.co.jp/mpmeister/>

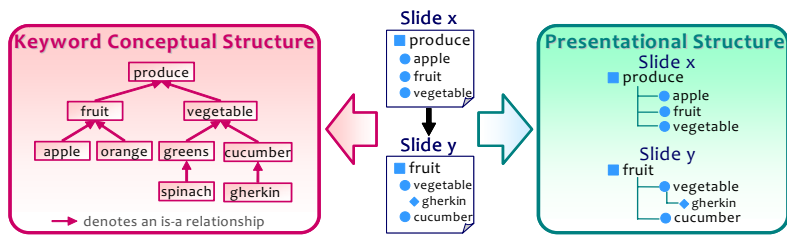


Fig. 1. Example of a relation between slides

Moreover, one of the important functions necessary for archiving presentation slides is to be able to retrieve desired slides, which are the important slides retrieved by the given keywords. In fact, for the benefit of the users, it is very important that certain keywords are supported, which will help them retrieve important slides. However, only retrieving the important slides on the basis of certain keywords can destroy the relevant information implicitly between slides will trend to lose the context that cannot help users' understanding. We need to retrieve appropriate slides of users desiring to learn some concepts represented by the query, easily.

Retrieving slides to meet users' requirements can be mainly achieved by (1) understanding the relation between slides in terms of a user query and (2) ranking the retrieved slides on the basis of the relation between slides related to the query. We find that a semantic relationship, e.g., an is-a, a part-of relationship, implicitly exists between keywords extracted from a slide text. Then, we derive a keyword conceptual structure by using the conceptual relationship existing between keywords extracted from the slide text. On the other hand, the usage of keywords in the slide differs depending on the author. Then, we derive a presentational structure by focusing on certain features of the slides, such as the level position information of indents in the slide text. Thus, it would be necessary to use the conceptual relationship and presentational structure to determine the relations between slides in terms of keywords.

Some presentation slides may be related to other slides in terms of detailed relation or generalized relation. For example, the explanation provided in slide y on "cucumber and gherkin of the cucurbitaceous vegetable" is more likely to be a detailed one than a general one provided in slide x on "vegetable". Therefore, we considered that the relation between slide x and slide y is a *detailed* relation in terms of "vegetable" (see Fig. 1).

Thus, we analyze the relationship between keywords and how the keywords vary in different level position of indents in the slide related to a query. We define the keyword conceptual structure of keywords extracted from the slide text by using the conceptual dictionary WordNet [1] is focused on the conceptual relationship, e.g., an is-a relationship, and we define the presentational structure of indents in the slide text. In addition, we provide two measures for ranking slides on the basis of the relations between slides related to a query. As mentioned above, we believed that an efficient presentation content retrieval engine would retrieve the slides that provided relevant information about the query and rank the retrieved slides on the basis of special measures.

The next section reviews related works. Section 3 explains the keyword conceptual structure and presentational structure, and mathematically determines the relations between slides. Section 4 describes the ranking measures for retrieving slides by using the relations between slides, and Section 5 explains and evaluates our prototype system. Finally, Section 6 presents the conclusions of this paper.

## 2. Related Works

Most of the research related to academic contents has focused on retrieval of slides. Yokota et al. [2] proposed a system termed Unified Presentation Slide Retrieval by Impression Search Engine (UPRISE) for retrieving a sequence of desired slides from archives containing a combination of slides and recorded videos. Kobayashi et al. [3] proposed a method of using laser pointer information for retrieving slides of a lecture by UPRISE. Le et al. [4] proposed a method for extracting important slides for automatically generating digests from the recorded presentation videos. Their method extracts important slides from unified content on the basis of the metadata features of a single medium or two heterogeneous media. However, we considered that only retrieving the important slides can lose the context in terms of user queries that cannot help users' understanding. Therefore, our objectives are to retrieve users' desired slides effectively by attracting the relevant information implicitly between slides in

terms of user queries, and to rank the retrieved slides into semantic orders by using the relations between slides related to user queries.

Kitayama et al. [5] proposed a method for extracting scenes on the basis of their relations and roles. Wang et al. [6] presented a method for automatically generating learning channels by using the semantic relations that implicitly exist in slides of a lecture that has an accompanying recorded video. These studies are similar to our study, where a method for retrieving desired slides using the relations between slides that attracting the relevant information between slides is proposed. However, we not only use semantic relations for slide retrieval but also focus on ranking the retrieved slides on the basis of the relations between slides.

Tanaka et al. [7] focused on the manipulation of complex objects and introduced the concept of “element-based” generalization relationships between complex objects and two new abstraction operators, namely, reduction and unification operators. Lan et al. [8] presented a theoretical framework for ranking the retrieved slides and demonstrated a method to perform generalization analysis of list-wise ranking algorithms using the framework. In our approach, however, we focused on generalization relationships between keywords in the slide text and utilized indents in the slide containing the keywords. Then, we determined the relations between slides and the keywords in order to retrieve the desired slides, and our ranking algorithm on the basis of measures that the retrieved slides into semantic orders by using the relevant information, i.e., whether the slides contained detailed information (DETAIL) or whether they contained general information (GENERALITY).

### 3. Determining the Relations between Slides using Conceptual Relationship and Presentational Structure

#### 3.1. Keyword Conceptual Structure and Presentational Structure

We consider that a semantic relationship implicitly exists between keywords extracted from the slide text. In particular, the conceptual relationship is called an is-a relationship [9, 10] and is used as a basis of the semantic relationship between keywords. “X subsumes Y, or Y is-subsumed-by X” (Y is-a X) usually means that concept Y is a specialization of concept X, and concept X is a generalization of concept Y. For example, a “fruit” is a generalization of an “apple”, an “orange”, a “mango”, and many other fruits, i.e., an apple is a fruit (apple is-a fruit). Therefore, we define a keyword conceptual structure consists of an is-a relationship between keywords are extracted using WordNet [1].

We define a presentational structure on the basis of indents in the slide text. The slide title (1st level indent) is the upper level. The first item of the text is on the 2nd level, and subitems deepen with the level of indentation (3rd level, 4th level, and so on). Indents outside the text, such as figures or tables, are on the average level of the slide. If a given keyword appears in the title of the slide or in less-indented lines, we implicitly assume that the lower-level indented keywords are supplementary and they explain the upper-level keywords.

Therefore, the conceptual relationship between keywords and the level position information of indents in a slide should be considered for slide retrieval.

#### 3.2. Determination of Relation Types

We define a focused slide and other slides that have specific relations as being conceptually related to the focused slide through one of the two types of relations: *detailed* and *generalized* relations. If a slide has a detailed relation with the other slides, we call this slide a detailed slide. Further, if a slide has a generalized relation with the other slides, we call this slide a generalized slide.

This section explains the manner in which the types of relations are determined. Let  $x$  be the number of a focused slide, and  $y$  be the number of a slide we want to retrieve. Slide  $x$  contains keywords  $k_i$  and  $k_m$ . The types of relations are determined for all slides for keyword  $q$  in a user query.

##### 3.2.1. Determination of Detailed Relations

If a slide has more information about a user query than the focused slide, its relation with the focused slide is a *detailed* one. We explain the determination of *detailed* slides by using the query keyword  $q$  present in the focused slide  $x$  and the slide  $y$ , which needs to be retrieved. Fig. 2 shows an example of determining the *detailed* relations between slide  $x$  and slide  $y$  for a query on a “vegetable”.

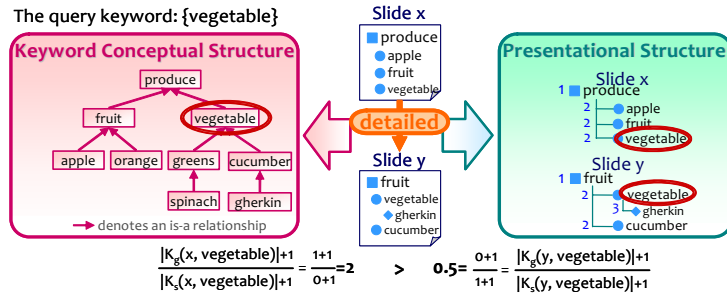


Fig. 2. Example of detailed relation between slides

When the query keyword  $q$  and other keywords in slide  $x$  and slide  $y$  conform to certain conditions, slide  $y$  is determined to be the detailed slide of slide  $x$ . This is because  $q$  appears more frequently in slide  $y$  than it does in slide  $x$ .

$$K_g(x, q) = \{k_i \mid k_i \in x, \text{level}(q) \geq \text{level}(k_i), q \text{ is-a } k_i\} \tag{1}$$

$$K_s(x, q) = \{k_m \mid k_m \in x, \text{level}(q) < \text{level}(k_m), k_m \text{ is-a } q\} \tag{2}$$

Where  $K_g(x, q)$  is a set of keywords in the slide  $x$  such that their level positions are not lower than the level position of  $q$  in the presentational structure, and  $q$  is-a each one of them in the keyword conceptual structure. In Eq. (1),  $k_i$  belongs to the set of keywords  $K_g(x, q)$  in the slide  $x$  and that its level position is not lower than that of  $q$  in the presentational structure, and  $q$  is-a  $k_i$  in the keyword conceptual structure. In our method, we extract the keyword conceptual structure as a tree-shaped structure. In general, an is-a relationship between keywords is equivalent to a parent-child relationship, and our method is susceptible to an is-a relationship as a descendant relationship.  $K_s(x, q)$  is a set of keywords in slide  $x$ , and their level positions are lower than the level position of  $q$  in the presentational structure, and each keyword is-a  $q$  in the keyword conceptual structure. In Eq. (2),  $k_m$  belongs to the set of keywords  $K_s(x, q)$  in slide  $x$  and that its level position is lower than that of  $q$  in the presentational structure, and  $k_m$  is-a  $q$  in the keyword conceptual structure.

$$\frac{|K_g(x, q)|+1}{|K_s(x, q)|+1} > \frac{|K_g(y, q)|+1}{|K_s(y, q)|+1} \tag{3}$$

Where the function  $|K_g(x, q)|$  extracts the total number of  $k_i$  in  $K_g(x, q)$ , and  $|K_s(x, q)|$  extracts the total number of  $k_m$  in  $K_s(x, q)$  in slide  $x$ .  $K_g(y, q)$  is also a set of keywords in slide  $y$ , satisfying the same conditions of  $K_g(x, q)$  by Eq. (1), and  $K_s(y, q)$  is a set of keywords in slide  $y$ , satisfying the same conditions of  $K_s(x, q)$  by Eq. (2). Thus, Eq. (3) can be used to calculate the ratio of  $|K_g(x, q)|$  to  $|K_s(x, q)|$  for slide  $x$  and the ratio of  $|K_g(y, q)|$  to  $|K_s(y, q)|$  for slide  $y$ . If the ratio calculated for slide  $x$  is higher than that calculated for slide  $y$  by Eq. (3), slide  $y$  is determined to be the detailed slide of slide  $x$  with regard to  $q$ .

### 3.2.2. Determination of Generalized Relation

If a slide contains content about the query in the outline given in a *generalized* slide, it is described in relation to the focused slide. We explain the determination of *generalized* slides by using the query keyword  $q$  present in the basic slide  $x$  and slide  $y$ , which needs to be retrieved.

$$\frac{|K_g(x, q)|+1}{|K_s(x, q)|+1} < \frac{|K_g(y, q)|+1}{|K_s(y, q)|+1} \tag{4}$$

When the query keyword  $q$  and other keywords in slide  $x$  and slide  $y$  conform to Eqs. (1), (2), and (4), then slide  $y$  is determined to be a generalized slide of slide  $x$ . This is because  $q$  appears more frequently in slide  $y$  than it does in slide  $x$ . Eq. (4) can be used to calculate the ratio of  $|K_g(x, q)|$  to  $|K_s(x, q)|$  for slide  $x$  and the ratio of  $|K_g(y, q)|$  to  $|K_s(y, q)|$  for slide  $y$ . When the ratio calculated for slide  $x$  is lower than that calculated for slide  $x$  by Eq. (4), slide  $y$  is determined to be the generalized slide of slide  $x$  with regard to  $q$ .

As can be seen, the detailed and generalized slides are functionally interchangeable, whereas a focused slide is a generalized slide from the viewpoint of a detailed slide.

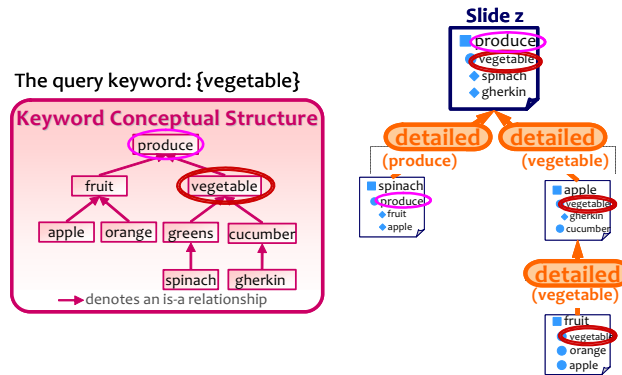


Fig. 3. Example of a slide with a high degree of DETAIL

#### 4. Ranking of Retrieved Slides

Our proposed method retrieves slides by determining the relations between slides about a user query. It is also difficult for users to understand relevant information between the retrieved slides in terms of the user query. Moreover, we consider the users in the different levels of understanding that they have different desires on the retrieved slides in terms of the user query. In this paper, our method provides two types of semantic rankings that focus on two measures, namely, *DETAIL* and *GENERALITY*. This section describes how to calculate the degrees of these two measures for ranking slides on the basis of the relations between slides related to the user query by using the keyword conceptual structure and presentational structure.

##### 4.1. Slide Ranking based on the Measure of DETAIL

In an order of retrieved slides providing detailed information in terms of a user query, the user must have a deep understanding of the desired slides related to the query. Then, slide ranking by using the measure of *DETAIL* can aid the user to understand the query with a detailed explanation well. If a slide provides detailed information regarding a query as compared to that provided by other slides, this slide is known as a specific slide, and it provides specific explanation about other slides with a high degree of *DETAIL*. As shown in Fig. 3, slide  $z$  has a detailed relation with other slides in terms of the content on “produce” and “vegetable” related to the query keyword “vegetable”.

We consider the function of the degree of *DETAIL* using the following indicators.

- The number of the target slide  $x$  has a detailed relation with the generalized slide  $G(x,q)$  with regard to the query keyword  $q$ .
- The generalized keyword  $k_c$  of  $q$  (means  $q$  is-a  $k_c$ ) in  $x$  is extracted from the keyword conceptual structure.
- The relevance of  $k_c$  and  $q$  is expressed in terms of the distance between the position of  $k_c$  and  $q$  in the keyword conceptual structure.
- The distance between  $z$  and  $G(x,q)$  for  $q$  indicates the number of detailed relations existing between  $x$  and  $G(x,q)$ ; the distance between  $z$  and  $G(x,k_c)$  for  $k_c$  indicates the number of detailed relations existing between  $x$  and  $G(x,k_c)$ .

We described these indicators as follows. If  $x$  has a detailed relation with  $G(x,k_c)$  with regard to  $k_c$ , we can say that  $x$  includes detailed specific information regarding  $q$ . If the distance between the position of  $k_c$  and  $q$  is short in the keyword conceptual structure, then the relevance of  $k_c$  and  $q$  is high such that the value of relevance of  $x$  and  $G(x,k_c)$  is high. Further, if the number of  $G(x,q)$  and  $G(x,k_c)$  is large and the number of detailed relations between  $x$  and  $G(x,q)$  with regard to  $q$  or that between  $x$  and  $G(x,k_c)$  with regard to  $k_c$  is large, then the distance between them is long such that the value of *DETAIL* of  $x$  is high.

The function of the degree of *DETAIL* is expressed as

$$det\_Val(x) = \sum_{q \in x} G(x,q) \times dist(x,G(x,q)) + \sum_{k_c \in x, q \text{ is-a } k_c} \frac{G(x,k_c) \times dist(x,G(x,k_c))}{pos(q) - pos(k_c) + 1} \quad (5)$$

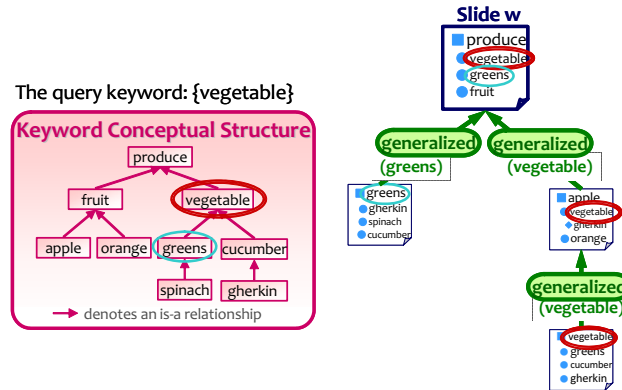


Fig. 4. Example of a slide with a high degree of GENERALITY

Where  $G(x,q)$  that extracts the number of  $x$  has a detailed relation of  $G(x,q)$  with regard to  $q$ . Further,  $dist(x, G(x,q))$  is the distance between  $x$  and  $G(x,q)$  with regard to  $q$  that extracts the number of detailed relations between  $x$  and  $G(x,q)$ . It should be noted that  $G(x,k_c)$  that extracts the number of  $x$  has a detailed relation of  $G(x,k_c)$  about  $k_c$  with regard to  $q$ . The function  $pos(q)-pos(k_c)+1$  is the relevance of  $k_c$  and  $q$  that extracts the distance between the position of  $k_c$  and  $q$  in the keyword conceptual structure. Further, the function  $dist(x,G(x,k_c))$  extracts the distance between  $x$  and  $G(x,k_c)$  in terms of the number of detailed relations between  $x$  and  $G(x,k_c)$ . Then, the function calculates the relevance of  $x$  and  $G(x,q)$  with regard to  $q$ , and the relevance of  $x$  and  $G(x,k_c)$  for  $k_c$  with regard to  $q$ . If  $G(x,q)$  and  $dist(x,G(x,q))$  are large, the value of detail between  $x$  and  $G(x,q)$  is high. If  $pos(q)-pos(k_c)+1$  is small, the value of relevance of  $k_c$  and  $q$  is high. If  $G(x,k_c)$  and  $dist(x,G(x,k_c))$  are large, the degree of *DETAIL* of  $x$  is high.

#### 4.2. Slide Ranking based on the Measure of GENERALITY

In an order of retrieved slides providing general information in terms of a user query, the user must easily grasping the general-content of the desired slides related to the query. Then, slide ranking on the basis of the measure of *GENERALITY* can aid the user to obtain a generalized explanation about the query, easily. If a slide provides general information regarding a query as compared to that provided by other slides, this slide is known as a general slide, and it provides explanation about other slides with a high degree of *GENERALITY*. As shown in Fig. 4, slide  $w$  has a generalized relation with other slides in terms of the content on “vegetable” and “greens” related to the query keyword “vegetable”.

We consider the function of the degree of *GENERALITY* using the following indicators.

- The number of target slide  $x$  has a generalized relation with the detailed slide  $D(x,q)$  with regard to the query  $q$ .
- The specified keyword  $k_p$  of  $q$  (means  $k_p$  is-a  $q$ ) in  $x$  is extracted from the keyword conceptual structure.
- The relevance of  $k_p$  and  $q$  is expressed in terms of the distance between the position of  $k_p$  and  $q$  in the keyword conceptual structure.
- The distance between  $x$  and  $D(x,q)$  for  $q$  indicates the number of generalized relations existing between  $x$  and  $D(x,q)$ ; the distance between  $x$  and  $D(x,k_p)$  for  $k_p$  indicates the number of generalized relations existing between  $x$  and  $D(x,k_p)$ .

We described these indicators as follows. If  $x$  has a generalized relation of  $D(x,k_p)$  with regard to  $k_p$ , we can say that  $x$  includes general information regarding  $q$ . If the distance between  $k_p$  and  $q$  is short in the keyword conceptual structure, then the relevance of  $k_p$  and  $q$  is high such that the value of relevance of  $x$  and  $D(x,k_p)$  is high. Further, if the number of  $D(x,q)$  and  $D(x,k_p)$  is large and the number of generalized relations between  $x$  and  $D(x,q)$  with regard to  $q$  or that between  $x$  and  $D(x,k_p)$  with regard to  $k_p$  is large, then the distance between them is long such that the value of *GENERALITY* of  $x$  is high.

The function of the degree of *GENERALITY* is expressed as

$$gen\_Val(x) = \sum_{q \in x} D(x,q) \times dist(x, D(x,q)) + \sum_{k_p \in x, k_p \text{ is-a } q} \frac{D(x,k_p) \times dist(x, D(x,k_p))}{pos(k_p) - pos(q) + 1} \quad (6)$$

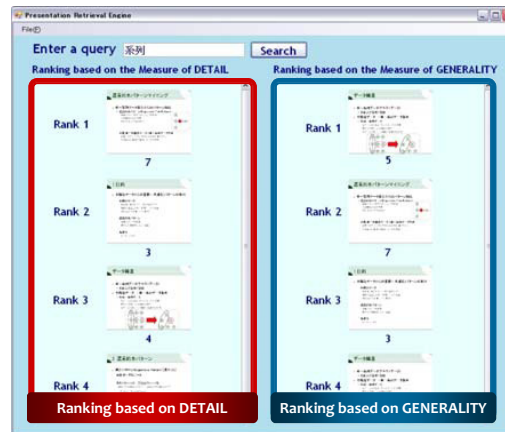


Fig. 5. Screenshot of prototype system

Where  $D(x,q)$  that extracts the number of  $x$  has a generalized relation of  $D(x,q)$  with regard to  $q$ . Further,  $dist(x,D(x,q))$  is the distance between  $x$  and  $D(x,q)$  with regard to  $q$  that extracts the number of generalized relations between  $x$  and  $D(x,q)$ . It should be noted that  $D(x,k_p)$  that extracts the number of  $x$  has a generalized relation of  $D(x,k_p)$  about  $k_p$  with regard to  $q$ . The function  $pos(k_p)-pos(q)+1$  is the relevance of  $k_p$  and  $q$  that extracts the distance between the position of  $k_p$  and  $q$  in the keyword conceptual structure. Further, the function  $dist(x,D(x,k_p))$  extracts the distance between  $x$  and  $D(x,k_p)$  in terms of the number of generalized relations between  $x$  and  $D(x,k_p)$ . Then, the function calculates the relevance of  $x$  and  $D(x,q)$  with regard to  $q$ , and the function calculates the relevance of  $x$  and  $D(x,k_p)$  about  $k_p$  with regard to  $q$ . If  $D(x,q)$  and  $dist(x,D(x,q))$  are large, the value of generality between  $x$  and  $D(x,q)$  is high. If  $pos(k_p)-pos(q)+1$  is small, the value of the relevance of  $k_p$  and  $q$  is high. If  $D(x,k_p)$  and  $dist(x,D(x,k_p))$  are large, the degree of *GENERALITY* of  $x$  is high.

As can be seen, our method for retrieving users' desired slides and ranking the retrieved slides into two types focuses on different measures and can satisfy users' demands.

## 5. Evaluation

### 5.1. Prototype System

We developed a prototype system to support a presentation content retrieval engine (see Fig. 5) in Microsoft Visual Studio 2008 C#. This prototype implements the determination part and the output part. In the determination part, all types of relations between slides are determined on the basis of the keyword conceptual structure by using WordNet [1] extracts is-a relationship between keywords and the presentational structure by using the level position information about keywords in the slide. Further, slide rankings are determined by calculating the two degrees of measures on the basis of the relations between slides related to a user query. The terms in the slides are extracted using a morphological analyzer *Mecab*<sup>3</sup>, which is in SlothLib<sup>4</sup> [11]. Using this system, a user can select the presentation content for studying. When the user enters a query of interest in the textbox and presses the "Search" button, the retrieved slides in two ranking types are presented in the output part.

### 5.2. Experiment 1: Validity of Determining Relation Types

There were five participants freely described the relations existing between two slides, and they assessed 155 sets

<sup>3</sup> <http://mecab.sourceforge.net/>

<sup>4</sup> <http://www.dl.kuis.kyoto-u.ac.jp/slothlib/>

Table 1. Experimental results

		Results determined using our system		
		<i>detailed relation</i>	<i>generalized relation</i>	<i>other</i>
Correct answers	<i>detailed relation</i>	<b>74</b>	9	35
	<i>generalized relation</i>	28	<b>34</b>	25
	<i>other</i>	5	5	44

Table 2. Results of relations between slides

	<i>detailed relation</i>	<i>generalized relation</i>	all
Precision	69.2% (74/107)	70.8% (34/48)	69.7% (108/155)
Recall	58.3% (74/127)	37.2% (34/86)	50.7% (108/213)
F-measure	0.63	0.51	0.59

of containing any keywords sampled at random from 4 actual academic contents<sup>5</sup>. Table 1 lists the results of the classification. The vertical column shows the results obtained using the proposed system; the horizontal rows show the correct answers given by the participants. We evaluated the coverage calculated using the slide set, which was determined to be any slide type identified by our system. The *others* cannot be determined by our system. The coverage reached a low of 63.6% by using our method. We therefore considered that the relation between content in presentations could not be expressed comprehensively only by using our method. This experiment confirmed the relations between slides containing any keyword could be cover by using the concept of relation types. We should improve the definitions of *detailed* or *generalized*, because they include a “parallel” relation and an “instance” relation; however, they were not frequent that difficult to define. In our method, we focused on *detailed* or *generalized*, and we expanded it to determine other types as semantic relations.

In addition, we evaluated the validity of the rules for determining the two types of relations defined by our method by precision and recall using the results obtained with the system and the results obtained from participants who gave correct answers. The results of the relations between slides are listed in Table 2. The recall of *detailed* or *generalized* was low, and many correct answers are detected to no relations by our method. We consider that the limitation of WordNet [1] is one factor that may cause the recall to be low. Although WordNet [1] is a large lexical database, it does not necessarily contain all concepts about any experimental keyword. Further, in the case of *generalized*, even if the same slide set is considered, participants’ answers differ from being “generalized” and “parallel”. If participants answered “generalized”, it means that they can understand the content well. However, if the participants answered “parallel”, it means that they understudied that slides have a relation at a minimum. We consider *generalized* to be effective when a user can understand that the slides have a relation at a minimum, but cannot determine the relation types.

This experiment confirmed that slides in the academic content have some kinds of relations between each other. Our proposed relations might provide an appropriate definition of using the conceptual relationship between keywords and the presentational structure of indents. Further, we believe that a considerable number of slides in the academic content provide detailed explanations. However, we should use an enhanced method for extracting the conceptual relationship between keywords; this method should not involve the use of WordNet [1] only, such as involving the use of a large ontology construction from Wikipedia.

### 5.3. Experiment 2: Validity of Ranking Types

We showed the participants the following ten data sets that are ten query keywords in the retrieved slides from 4

<sup>5</sup> DBSJ Archives: <http://www.dbsj.org/Japanese/Archives/archivesIndex.html>



Table 3. Comparison between the Spearman's rank correlation coefficients obtained by participant evaluation and our method

Data set	Query keyword	Ranking type	Participants						Average
			A	B	C	D	E	F	
(1)	“pattern”	<i>DETAIL</i>	0.77	0.63	0.80	0.80	0.63	0.77	0.73
(2)	“order”	<i>DETAIL</i>	1.00	0.70	0.70	0.70	0.70	0.6	0.73
(3)	“data”	<i>DETAIL</i>	1.00	0.40	0.40	0.20	0.80	0.80	0.60
(4)	“composition”	<i>DETAIL</i>	0.80	0.20	1.00	0.80	0.40	0.20	0.57
(5)	“contrast”	<i>DETAIL</i>	1.00	0.50	1.00	1.00	0.50	1.00	0.83
(6)	“relationship”	<i>GENERALITY</i>	0.80	0.80	0.20	-0.40	0.60	0.40	0.40
(7)	“aggregate”	<i>GENERALITY</i>	0.80	0.80	0.80	1.00	0.20	0.40	0.67
(8)	“group”	<i>GENERALITY</i>	0.20	0.80	1.00	1.00	-0.20	1.00	0.63
(9)	“geography”	<i>GENERALITY</i>	1.00	0.40	-0.40	0.40	0.80	0.40	0.43
(10)	“comparison”	<i>GENERALITY</i>	0.70	-0.10	0.70	0.70	0.90	0.60	0.58

actual academic contents used in Experiment 1. Then, we let them rank the slides with regard to ten query keywords in the order of degree of *DETAIL* and *GENERALITY*, respectively. For each query keyword, we then calculated the Spearman's rank correlation coefficient between the participant rankings and our method rankings. The Spearman's rank correlation coefficient ranges from -1 to 1, where -1 indicates that two rankings are completely reverse whereas 1 indicates that the rankings are exactly the same.

The experimental results are listed in Table 3. Six participants, i.e., A to F, participated in this experiment. We can see that the degrees of our proposed measures, i.e., *DETAIL* and *GENERALITY*, are greater than 0 and that on an average, the measure determined by the participants, i.e., *DETAIL*, shows the best performance. These results indicate that our proposed method that takes into account the relation between slides about the query keyword based on the conceptual relationship between keywords and the presentational structure of indents can be successfully applied to the presentation content retrieval engine on the basis of semantic rankings. However, the degree of *GENERALITY* on an average was low here, i.e., data set (6) and data set (9). We calculated the degree of *GENERALITY* by using the *generalized* relation between slides with regard to the query keyword and the specified keyword of the query keyword. In particular, the degree determined by participant D in data set (6) and participant C in data set (9) were too low. We consider that it is difficult for participant D and participant C to ascertain the specified keywords of the query keyword in the retrieved slides, which may reduce the performance. Although a slide has a *generalized* relation with other retrieved slides with regard to the query keyword and it contains many specified keywords of the query keyword that has a *generalized* relation with other slides that were not retrieved, the specified keywords of the query keyword were unknown by a participant. It can be seen that our method can extract many concepts of the query by effectively using the conceptual relationships between keywords. From the results of this experiment, we find that we have to improve the determination of the ranking algorithm by using the relations between slides containing keywords and the conceptual relationship between keywords.

## 6. Concluding Remarks

We have proposed a presentation content retrieval engine that uses the conceptual relationship and presentational structure. The slide ranking method is used to retrieve slides and rank the retrieved slides on the basis of the relations between slides with regard to a user query. The type of relation is determined on the basis of the conceptual relationship between keywords and the presentational structure in the slide text of keywords. Thus, users use the presentation content retrieval engine to retrieve desired slides by a user query and rank the retrieved slides on the basis of two measures. Moreover, we have also developed a prototype system and evaluated it using actual presentation data. We confirmed an improvement in the coverage of the types of relations and their definitions and the validity of slide rankings by using the relations between slides related to the query keywords.

In future, we intend to improve the algorithm to determine the relations between slides in order to calculate the measures of slide rankings. Furthermore, our method can enhance the retrieval technique that would be useful for a

user if he or she asks a query that includes two and more keywords; we have to determine the relationship between keywords in the query to retrieve the user's desired slides by analyzing the relevance of the queried keywords and develop other measures for ranking slides from multiple presentation contents. Therefore, we intend to evaluate the effect of our proposed indicators in the ranking method using a presentation content, and then compare its efficiency with the traditional  $tf \cdot idf$  retrieval method using multiple presentation contents. Moreover, we plan to evaluate the effectiveness of our method with a large set of actual lecture contents.

We can also extend our approach application areas, such as Web retrieval. Keyword Conceptual structure and presentational structure in the text document cannot directly be used for searching Web pages, because they are features peculiar to presentation contents. However, we can consider possibility of the application of some parts of our approach by utilizing the information about the conceptual relationship between keywords and the hierarchical structure in the text document of browsing Web pages. With in regard to ranking, Google also uses position information to rank retrieved Web pages [12], with font and capitalization information. We currently do not use font nor capitalization information, but it is not difficult to improve our approach with these information. It is important to consider the abstraction of presentation content to make retrieval more intelligent.

## Acknowledgements

This research was supported in part by a Grant-in-Aid for Scientific Research (B)(2) 20300039 from the Ministry of Education, Culture, Sports, Science, and Technology of Japan.

## References

1. WordNet, <http://wordnet.princeton.edu/>.
2. H. Yokota, T. Kobayashi, H. Okamoto, and W. Nakano, "Unified contents retrieval from an academic repository," in *Proc. of International Symposium on Large-scale Knowledge Resources (LKR2006)*, pp. 41–46, March 2006.
3. T. Kobayashi, W. Nakano, H. Yokota, K. Shinoda, and S. Furui, "Presentation scene retrieval exploiting features in videos including pointing and speech information," in *Proc. of International Symposium on Large-scale Knowledge Resources (LKR2007)*, pp. 95–100, March 2007.
4. H.-H. Le, T. Lertrudachakul, T. Watanabe, and H. Yokoda, "Automatic digest generation by extracting important scenes from the content of presentations," in *Proc. of the 19th International Conference on Database and Expert Systems Application (DEXA2008)*, pp. 590–594, September 2008.
5. D. Kitayama, A. Otani, and K. Sumiya, "A scene extracting method based on structural and semantic analysis of presentation content archives," in *Proc. of the 7th International Conference on Creating, Connecting and Collaborating through Computing (C52009)*, pp. 275–280, January 2009.
6. Y. Wang, D. Kitayama, R. Lee, and K. Sumiya, "Automatic generation of learning channels by using semantic relations among lecture slides and recorded videos for self-learning systems," in *Proc. of the 11th IEEE International Symposium on Multimedia (ISM2009)*, pp. 275–280, December 2009.
7. K. Tanaka and M. Yoshikawa, "Towards abstracting complex database objects: Generalization, reduction and unification of set-type objects (extended abstract)," *Proc. 2nd Int. Conf. on Database Theory (Lecture Notes in Computer Science)*, vol. 326, pp. 252–266, August 1988.
8. Y. Lan, T. Liu, Z. Ma, and H. Li, "Listwise approach to learning to rank: Theorem and algorithm," in *Proc. of the 26th Annual International Conference on Machine Learning (ICML2009)*, pp. 577–584, June 2009.
9. G. A. Miller, R. Beckwith, C. Fellbaum, D. Gross, and K. J. Miller, "Introduction to wordnet: An on-line lexical database," in *International Journal of Lexicography*, pp. 235–244, 1999.
10. G. A. Miller, "Wordnet: A lexical database for english," *Communications of the ACM (CACM)*, vol. 38, no. 11, pp. 39–41, November 1995.
11. H. Ohshima, S. Nakamura, and K. Tanaka, "Slotlib: A programming library for research on web search," *Proc. of the Database Society of Japan (DBSJ Letters)*, vol. 6, no. 1, pp. 13–116, June 2007 (in Japanese).
12. S. Brin and L. Page, "The autonomy of a large-scale hypertextual web search engine," *Proc. of the 7th WWW Conf.*, vol. 30, pp. 107–117, April 1998.