

Construction of Personalized and Purpose-Oriented Thesaurus



Masaharu YOSHIOKA

Hokkaido Univ.

Makoto HARAGUCHI



Background and Objectives

■ Background:

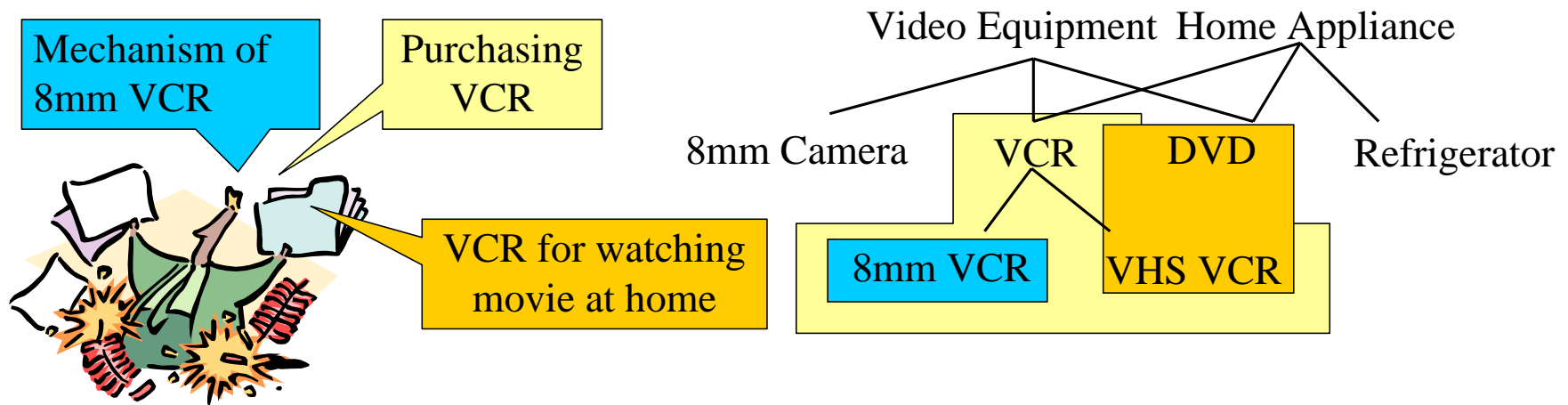
- Concept structure of electronic dictionary and thesauri
 - General purpose thesaurus is too general and they do not have fine granularity abstract concepts to represent users' particular intention
 - Concept category for representing user's particular intention is selective and is not directly corresponds to a single common concept category stored in general thesauri

■ Objectives

- Propose methodology to construct personalized and purpose-oriented thesaurus
 - Modification of concept structure of electronic dictionary for particular user's intention
 - Evaluation of Construct Personalized and Purpose-Oriented Thesaurus by Application to Information Retrieval System

Appropriate Concept Category for Particular Users Intention

- Role of keyword in different situation
 - Use a keyword as a particular concept
 - A keyword corresponds to a concrete concept category
 - Use a keyword as a concept category
 - A keyword corresponds to an abstract concept category
 - Use a keyword as representative concept of intended category
 - A keyword corresponds to an abstract concept category
 - General thesaurus may not have appropriate concept category



Construction of Tentative Concept Category

- Construction of concept category from relevant documents
 - Assumption
 - Tentative concept category should exist more frequently than usual documents
 - Alternative keyword for representing this category exists in relevant documents
 - Method
 - Use mutual information content for selecting an abstract concept category
 - Collect keywords that belongs to selected category to make tentative concept category
 - Evaluation measure for selecting an abstract concept category
 - Appropriate concept category for characterizing relevant documents might have higher mutual information content with relevant documents

$$I(T;W) = I(W;T) = H(W) - H(W|T)$$
$$= \sum_{w \in W} \left[-p(w) \log_2 p(w) + \sum_{t \in T} p(t) p(w|t) \log_2 p(w|t) \right] \dots (1)$$

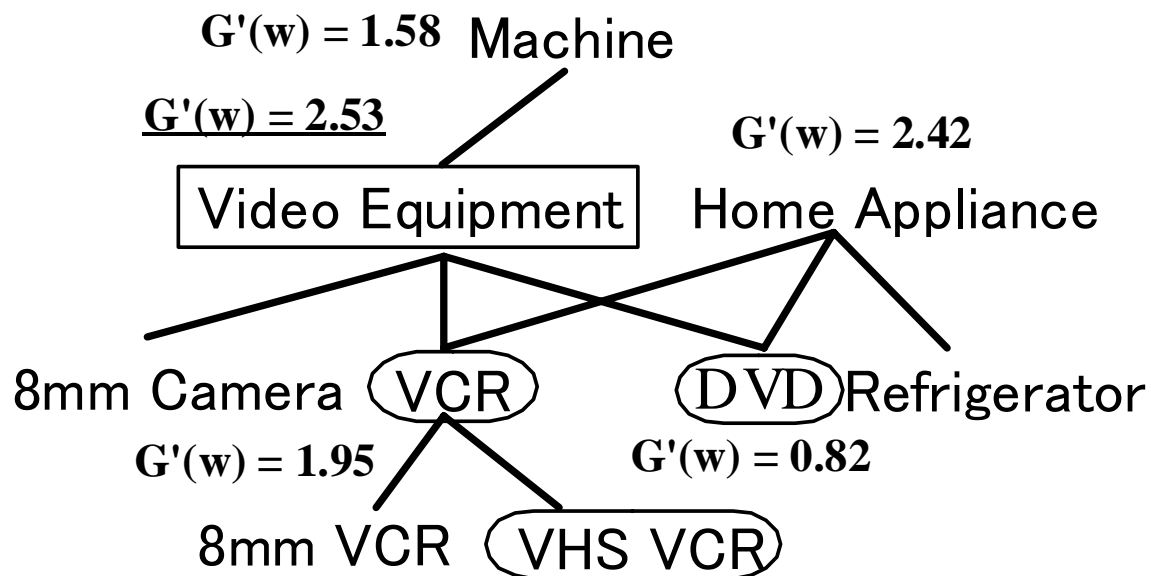
$$G(w) = -p(w) \log_2 p(w) + \sum_{t \in T} p(t) p(w|t) \log_2 p(w|t)$$
$$= p(r) p(w|r) \{ \log_2 p(w|r) - \log_2 p(w) \}$$
$$+ p(\bar{r}) p(w|\bar{r}) \{ \log_2 p(w|\bar{r}) - \log_2 p(w) \}$$
$$\approx p(r) p(w|r) \log_2 \frac{p(w|r)}{p(w)} \dots (2)$$

Example of Tentative Concept Category

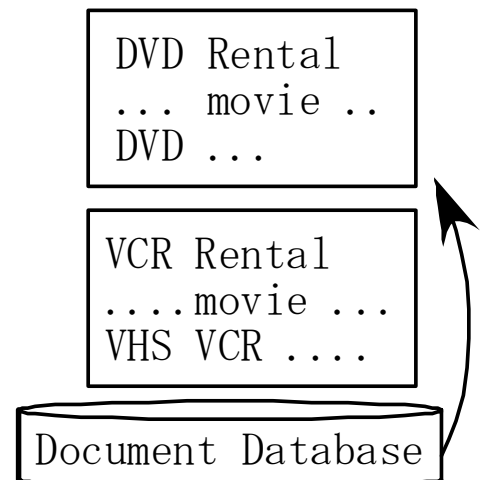
- Evaluation measure for selecting an appropriate concept category

$$G'(w) = p(w|r) \log_2 \frac{p(w|r)}{p(w)} \dots (3)$$

- Construction method
 - Calculate evaluation measure for each keyword and candidate abstract concept category
 - Select concept that has highest evaluation measure for appropriate one



Query: VCR for movie





IR System based on Personalized and Purpose-Oriented Thesaurus

- Evaluate effectiveness of tentative concept construction
- IR system that uses personalized and purpose-oriented thesaurus for query expansion
 - Prototype system
 - Implementation based on Information Retrieval Package (Uchiyama, 2001): Baseline system
 - BM25-base IR system
 - Electronic Dictionary
 - EDR Dictionary
 - Goi-Taikei- A Japanese Lexicon 日本語語彙体系
 - Test Collection: NTCIR-1 (Japanese)
 - Document set: Abstract for conference proceedings: 332918 documents
 - Number of topic:53
 - Type of data: Title, Description, Narrative, Concept
 - Relevance judgment: A (Relevant), B (Partially relevant), C (Irrelevant)
 - Type of retrieval:
 - long (Title+Description+ Narrative+Concept)
 - short(Description)
 - Use B relevance judgment and calculate mean-average precision

Prototype System

■ Construction of retrieval index

1. Apply morphological analyzer to extract initial index keywords (mostly from noun)
2. Check EDR (or Goi-Taikei) Dictionary about initial keywords and add concept ID list that corresponds to initial keywords to index

■ Information retrieval based on adaptive generalization

1. Extract initial query keywords by using morphological analyzer
2. By using relevant document sets (select 5 relevant documents by automatic feedback), the system construct tentative concept category
3. Add words of constructed category for query expansion

Initial Query

ビデオの紹介
ビデオは...

Chasen+
Baseline system

ビデオ 紹介
ビデオ ...

EDR

ビデオ 3c2a83 3c52fb ..
紹介 1faf20 3d04c8 ..

ビデオ 3c2a83
3c52fb ... 紹介
1faf20 3d04c8
... ビデオ 3c2a83
3c52fb ...

Final Index

Example of Constructed Concept Category

- Query title: 次世代インターネット(Next Generation Internet)
 - Original Keyword: 通信する(to communicate)
 - Generalized concept: 3cf631 23 subcategories and 56 words
 - Select 3 words from relevant documents

Query title: 次世代インターネット(Next Generation Internet)
Original Keyword: 通信する(to communicate)

3cf631 通信する(to communicate)



Tentative Category

通信する(to communicate)
放送する(to transmit a message to receivers through a wire or radio system)
中継する(to relay messages, speeches or music from a broad casting station)

0e64a6 衛星中継する(the act of relaying electric waves by artificial satellite)

0f07c6 交信する(to exchange communication with some one)

0f021b5 混信する(to get jammed)



3c1208 電話する(to call a person on the telephone)

3c54cd 放送する(to transmit a message to receivers through a wire ore radio system)

3d1b03 中継する(to relay messages, speeches or music from a broad casting station)

Evaluation Result (EDR)

- PP-oriented thesaurus results has better performance than original thesaurus
- PP-oriented thesaurus has poorer performance than baseline system

Mean Average Precision

System type	Long	Short
PP-oriented Thesaurus	0.4587	0.3587
Original Thesaurus	0.4283	0.3131
Baseline System	0.4880	0.4097

Expanded Query Size

System type	Long	Short
PP-oriented Thesaurus	50.89	8.50
Original Thesaurus	46.22	7.44
Baseline System	209.90	186.47

Evaluation Result (Goi-Taikei)

- Similar to EDR case
- Goi-Taikei has poorer performance than EDR one

Mean Average Precision

System type	Long	Short
PP-oriented Thesaurus	0.4352	0.3106
Original Thesaurus	0.4282	0.3087
Baseline System	0.4880	0.4097

Expanded Query Size

System type	Long	Short
PP-oriented Thesaurus	67.08	9.53
Original Thesaurus	36.28	6.70
Baseline System	209.90	186.47



Discussion

- Coverage of thesauri
 - Since these thesauri don't have good coverage on technical terms, baseline system that adds all terms in relevant documents works well
 - We need a mechanism to expand queries by using words in relevant documents
- EDR dictionary has better performance than Goi-Taikei
 - Goi-Taikei has approximately 3000 concept categories and EDR has approximately 50000 concept categories.
 - Because of that, Goi-Taikei does not have appropriate abstract concept categories for particular purpose



Conclusion

- Propose methodology to construct personalized and purpose-oriented thesaurus
 - Propose methodology to select relevant term list from general abstract concept category.
 - It improve the performance of IR system compared with original thesaurus.
 - However, it is not sufficient for IR performance compared to automatic relevance feedback query expansion.
- Apply this method to different application domain.
 - Construct purpose-oriented thesaurus for a category in open directory.