

AN INTELLIGENT SYSTEM FOR INFORMATION RETRIEVAL OVER THE INTERNET THROUGH SPOKEN DIALOGUE

Hiroya Fujisaki¹, Hiroyuki Kameda², Sumio Ohno¹, Takuya Ito¹, Ken Tajima¹, and Kenji Abe¹

¹ Department of Applied Electronics, Science University of Tokyo
2641 Yamazaki, Noda, 278 Japan

² Tokyo Engineering University, 1404-1 Katakura, Hachioji, 192 Japan

ABSTRACT

For the purpose of coping with the affluence of information available over the Internet, an efficient, accurate and user-friendly system for information retrieval is mandatory. This paper presents an intelligent system based on the use of spoken dialogue as the main channel for user-system interface, use of key concepts, processing of unknown words, automatic acquisition of various kinds of knowledge for improving the performance, and agent technologies for system realization. Details of functions required for the agents are also described.

1. INTRODUCTION

With the rapid progress of computer technology and world-wide development of information networks, a vast amount of information is now being generated, published, and stored at a number of sites distributed all over the world. Such an affluence of information, however, is useless or may even become harmful unless one has a means for rapidly retrieving the information that is truly necessary and appropriate. In this respect, conventional systems for information retrieval are far from being satisfactory, and tend to collect irrelevant information as well as to miss relevant information. These situations can be ascribed, partly to the difficulty for the user to identify and express his/her intention precisely, and partly to the difficulty for the system to infer the user's intention correctly. These difficulties can be greatly reduced by introducing spoken dialogue between the user and the system.

While keyword search is suited for retrieving information from databases not necessarily designed by a common principle, both the accuracy and efficiency tend to be low because of polysemy and synonymy of keywords. These difficulties can be overcome by using 'key concepts' [1] rather than keywords. In this case, information retrieval is based,

not on the surface forms of keywords, but on their semantic contents intended by the user. Difficulties arising from polysemy of keywords can be solved by spoken dialogue if all the keywords are 'known' to the system (i.e., already registered in the lexicon of the system).

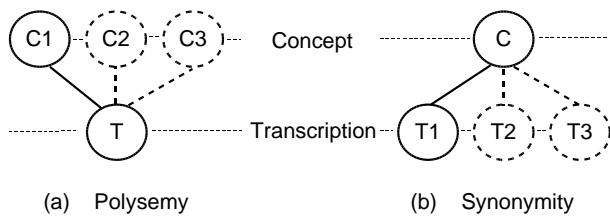
In actual information retrieval situations, where new words or new compound words made by combining known morphemes commonly occur, it is impossible that all the keywords are registered in the lexicon. Thus the system must have the ability to infer the meaning of new keywords that are 'unknown' (i.e., not registered in the lexicon). In other words, the system has to possess the ability of knowledge acquisition. This is also necessary if one aims at an intelligent system which will automatically improve its performance.

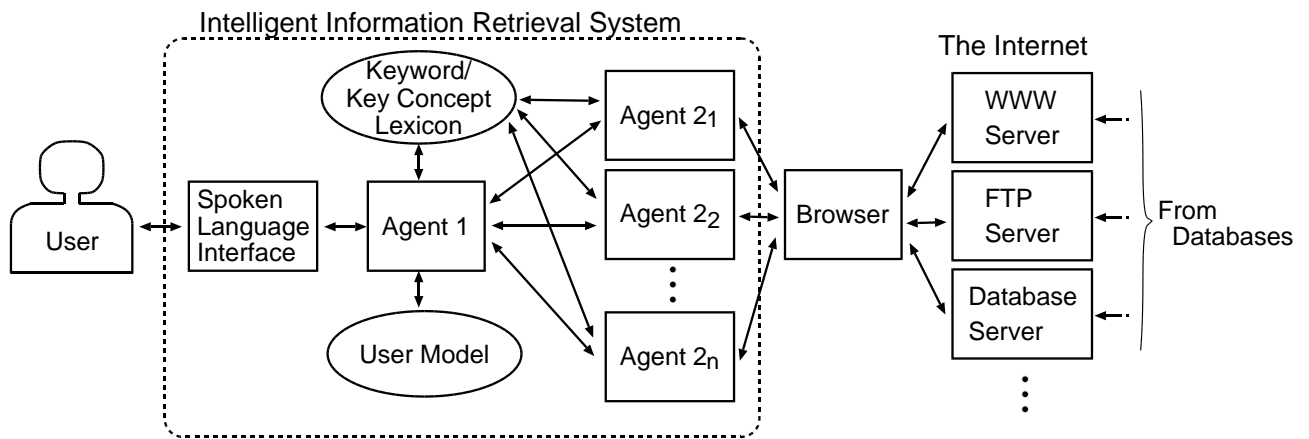
Based on these considerations, the present paper first describes the basic principles for advanced information retrieval, then proposes an intelligent system based on these principles, and describes some of the key functions of the system.

2. BASIC PRINCIPLES

2.1. Spoken Dialogue Between User and System

In many cases, a user is not fully aware, nor has sufficient knowledge, of the information which he/she wishes to retrieve. It is often the case that the user's intention becomes definite only after he/she gets some knowledge, through trial and error, from the system. If the user and the system can exchange information and knowledge through dialogue, especially through spoken language, it will greatly facilitate the process of formation and expression of intention on the part of the user, and also the process of its clarification on the part of the system.





3.2.6. Evaluation of Results Obtained by Agent 2

Using the weights mentioned in 3.2.1., Agent 1 evaluates the relevance of retrieved items, and those with higher relevance scores are presented first to the user. If they do not satisfy the user, those of lower relevance scores are shown. The final choice by the user is utilized to modify the weights and the model of the user for the purpose of improving performance. If the initial search does not satisfy the user, Agent 1 tries to obtain more key concepts through dialogue with the user. This process is repeated until the user is satisfied or the search is exhausted.

3.3. Functions of Agent 2

3.3.1. Information Retrieval from Database

Using the meta-knowledge on available databases, Agent 2 selects the databases to be accessed, and retrieves the requested information. The access is made through a web browser and relevant servers. Since the final stage of search is done through keywords, the search formula is translated into equivalent keywords, but the formula is far more precise and comprehensive than that would be obtained without going through key concepts.

3.3.2. Dealing with Ambiguities

Since the ambiguity due to polysemy of a keyword given in a data cannot be resolved through dialogue with the user, Agent 2 infers the most probable meaning on the basis of its context, i.e., in relation to other keywords / key concepts given in the data, and evaluates the relevance of the data.

3.3.3. Processing of Unknown Keywords in a Data

If a keyword given in a data is unknown, Agent 2 makes inference on its meaning by a procedure already developed [2]. It is based primarily on the structured analysis of the surface form, and the concept is inferred from those of the constituent morphemes.

3.3.4. Acquisition of Meta-Knowledge on the Databases

It is the task of Agent 2 to construct and maintain a meta-knowledge on the databases which it covers. Since the contents of many databases are frequently updated, Agent 2 conducts a regular check on each of the databases, and updates the meta-knowledge. Furthermore, it keeps record of the number of hits of each database as a meta-knowledge on its usefulness.

4. SUMMARY AND CONCLUSION

For the purpose of realizing a user-friendly and efficient system for information retrieval over the Internet, we have proposed a new system based on spoken dialogue as the main channel for user-system interface, use of key concepts, processing of unknown words, automatic acquisition of various kinds of knowledge for improving the performance, and agent technologies for system realization. Due to space limitations, many of the details have to be reported elsewhere. Work is under way to collect spoken dialogue data in various types of information retrieval problems, and to construct an efficient procedure for inferring the user's intention mainly through spoken dialogue, supplemented occasionally by a display and a keyboard.

ACKNOWLEDGMENT

The current work was supported by Japan Society for the Promotion of Science as a Project on 'Research for the Future' (Project No. JSPS-RFTF-96R15201).

REFERENCES

- [1] H. Fujisaki, H. Kameda and H. Kawai, "A system for information retrieval of newspaper articles based on key concepts," *Transactions on Natural Language Processing, Information Processing Society of Japan*, vol. 44, no. 4, 1984.
- [2] H. Kameda, H. Fujisaki, T. Morita and A. Kurashima, "Classification and processing of unknown words," *Proceedings of the 36th National Convention of the Information Processing Society of Japan*, pp. 1195–1196, 1988.
- [3] H. Fujisaki, et al., "An advanced information retrieval system based on extraction of key concepts and processing of unknown words," *Proceedings of the 54th National Convention of Information Processing Society of Japan*, vol. 3, pp. 23–24, 1997.
- [4] H. Fujisaki, et al., "An information retrieval system on the Internet using intelligent agents," *Proceedings of the 1997 General Conference of the Institute of Electronics, Information and Communication Engineers*, vol. 6, p. 261, 1997.
- [5] H. Fujisaki, et al., "Construction of a system for advanced information retrieval through spoken dialogue," *Proceedings of the Third Annual Meeting of the Association for Natural Language Processing*, pp. 261–264, 1997.