

Comparative Evaluation of a Natural Language Dialog Based System and a Menu Driven System for Information Access: a Case Study

Joyce Chai^{*}, Jimmy Lin⁺, Wlodek Zadrozny^{*}, Yiming Ye^{*}
Margo Budzikowska^{*}, Veronika Horvath^{*}, Nanda Kambhatla^{*},
Catherine Wolf^{*}

IBM T. J. Watson Research Center^{*} MIT Artificial Intelligence Lab⁺
30 Saw Mill River Rd. 545 Technology Square
Hawthorne, NY 10523 Cambridge, MA 02139
{jchai, wlodz, yiming, sm1, veronika, nanda, cwolf}@us.ibm.com jimmy@mit.edu

Abstract

This paper describes the evaluation of a natural language dialog based navigation system (HappyAssistant) that helps users access e-commerce sites to find relevant information about products and services. The prototype system leverages technologies in natural language processing and human computer interaction to create a faster and more intuitive way of interacting with websites, especially for the less experienced users. The result of a comparative study shows that users prefer the natural language enabled navigation two to one over the menu driven navigation. In addition, the study confirmed the efficiency of using natural language dialog in terms of the number of clicks and the amount of time required to obtain the relevant information. In the case study, comparing to the menu driven system, the average number of clicks used in the natural language system was reduced by 63.2% and the average time was reduced by 33.3%.

1 Introduction

With the emergence of e-commerce (Aggarwal & Wolf & Yu, 1998; Muller & Pischel, 1999), successful information access on e-commerce websites that accommodates both customer needs and business requirements becomes essential. Menu driven navigation and keyword search provided by most commercial sites have tremendous limitations. The menu driven approach is likely to overwhelm users and frustrate them with lengthy interactions. A recent study shows that the user's interest in a particular site decreases exponentially with the increase in the number of mouse clicks (Huberman & Pirolli & Pitkow, 1998). Therefore, shortening the interaction path to provide useful information becomes important. On the other hand, keyword search engines usually require users to know domain specific jargon. Keywords are not only unable to precisely describe the user's intention, but more importantly, they might not match words used in the catalog or documents. Furthermore, the keyword search lacks understanding of semantic meanings of the search words or phrases. For example, keyword search cannot understand that "summer dress" should be looked up in women's clothing under "dresses", whereas "dress shirt" most likely in men's under "shirts". A search for "shirt" can reveal dozens or even hundreds of items, which is useless for somebody who has a specific style and pattern in mind. Moreover, search engines do not accommodate business logic, e.g. a prohibition against displaying cheap earrings with more expensive ones. The solution to these problems lies, in our opinion, in centering electronic commerce websites on natural language (and multimodal) dialog. This claim is supported by results of a recent study we performed, and which will be presented in this paper.

We have built and evaluated a natural language (NL) dialog system for guiding users towards computer products: PCs, notebooks, servers and services, currently sold by IBM. The system allows customers to make requests in natural language and be directed towards appropriate web pages that sell the product or provide the service. Users can type in what they are looking for in natural language. The system identifies and understands key concepts from the user's input. Then by applying user's concepts to business rules, the system will either display the web page that satisfies user's requests or initiate a dialog with the user to either ask for additional information or clarify the request.

Natural language dialog has been used in many areas, such as for call-center/routing application (Carpenter & Chu-Carroll, 1998; Chu-Carroll & Carpenter, 1998), email routing (Walker & Fromer & Narayanan, 1998), information retrieval and database access (Androustopoulos & Ritchie, 1995), and for telephony banking (Zadrozny et al, 1998). The integration of natural language dialog with an e-commerce environment is a novel feature of our system. Our work goes beyond the "natural language interface" features of websites such as www.askjeeves.com and www.neuromedia.com, which work in a question-answer mode and do not use dialog. This is a crucial difference. When searching e-commerce sites, users often do not know where to find information, or how to specify a request although they have targets in their minds. Sometimes they only have vague or no targets in minds (Saito & Ohmura, 1998). Thus they need to formulate or revise their request based on additional information, which can be provided in a dialog. Our study shows that natural language dialog is a very effective medium for negotiating such contexts by understanding user's requests/intentions and providing help/advice/recommendations to the user.

Furthermore, information access on e-commerce sites is different from traditional keyword search or information retrieval where the business logic is not supported. However, for a business to successfully operate online, certain business rules should be enforced to facilitate the search. In our system, business rules are incorporated in the dialog management. We apply XML (Bray & Paoli & Speberg-McQueen, 1998; Radev et al, 1999) to represent and manage Domain Lexicon (concepts) and business rules.

We have carried out a user study to evaluate the natural language dialog based system, particularly, in comparison with a menu driven system. The result shows that users prefer the natural language dialog mode of interaction two to one over the menu driven interaction. The preference is stronger for less experienced internet users.

In this paper, we will first describe the natural language dialog based system (HappyAssistant). Then we will report the results from the comparative evaluation of this system and a menu driven system. Finally we will discuss what we have learned from the study and propose future work.

2 The HappyAssistant

2.1 System Architecture

The architecture of the system supports multimodal dialog. For this case study, the prototype system was implemented for textual input, with the possibility of browsing non-textual information. However, our architecture is designed to support inputs from different channels and modalities including keyboard input and output, speech input and output over a telephone, speech input and output over a microphone, mouse input, pointing device input, dataglove, etc. The system consists of three major modules: Presentation Manager (PM), Dialog Manager (DM) and Action Manager (AM). The presentation manager is responsible for separating content from the presentation mode. In the prototype system, the presentation manager employs a natural language parser to transform user's natural language query into a logical form, and sends the logical form to the dialog manager. It is also responsible for obtaining the system's response

from the dialog manager and presenting it to the user. The dialog manager is responsible for determining the specific action(s) requested by the user and filling the parameters (e.g. the attributes of the computers users are interested) of the identified action by way of a dialog with the user. The Knowledge Base for business rules specifies the translation from user requests to action plans for the action manager to satisfy the requests, for example, retrieving information about particular computer models from the product catalog. The architecture is shown in Figure 1.

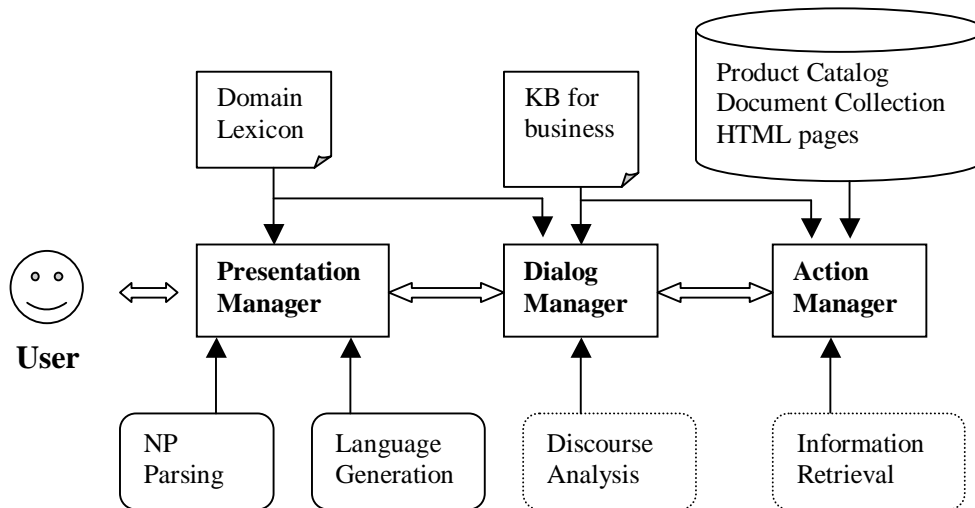


Figure 1: System Architecture

2.2 Domain Knowledge

We use XML to represent and manage the Domain Lexicon and the Knowledge Base for business rules. The Domain Lexicon is an external dictionary file that maps keywords onto concepts. Shown below is a fragment of that file:

```

<ENTRY NORMAL_FORM="affordable-concept"
  QUESTION="What are your financial constraints?"
  BACKOFF_QUESTION="Is affordable price important to you?">
  <WORD>affordable</WORD>
  <WORD>cheap</WORD>
  <WORD>inexpensive</WORD>
  <WORD>reasonably priced</WORD>
</ENTRY>
  
```

The <WORD> elements are keywords, which trigger the concept, and the QUESTION attribute is a natural language question designed to elicit that particular concept and related ones. Concepts can be as concrete as a category of computer, e.g. desktop-concept, or as abstract as an idea, e.g. performance-and-value.

A business rule consists of a list of concepts together with some metadata about the target product or service:

```

<RULE>
  <CONCEPT_LIST>
    <CONCEPT>notebook-concept</CONCEPT>
    <CONCEPT>high-tech-concept</CONCEPT>
    <CONCEPT>fast-concept</CONCEPT>
  </CONCEPT_LIST>
</RULE>
  
```

```

<WEIGHT>0.9</WEIGHT>
<DESCRIPTION>The IBM ThinkPad 770</DESCRIPTION>
<LONG_DESCRIPTION><![CDATA[
  <P><IMG SRC=http://...>
  <P>The ThinkPad 770 is IBM's top of the line laptop, offering the
  ultimate in performance and display.
]]></LONG_DESCRIPTION>
<URL><![CDATA[http://...]]></URL>
</RULE>

```

The <WEIGHT> element is designed to reflect how strongly the business wants to push this rule. In the case study, we manually set this value. However, in a real business setting, this weight should reflect customer demands, business decisions on pushing certain products, etc. For example, if the company wants to push ThinkPad 600 product, a business rule that translates a set of concepts to the ThinkPad 600 should have higher weight than business rules that translate the concepts to other products.

Knowledge management is a key issue in information systems. The knowledge base in our system should reflect the evolving business strategies and marketing messages. Although in the prototype system, the Domain Lexicon and the Knowledge Base for business rules were created manually, we are currently looking into developing tools to automate this process.

2.3 The Role of Natural Language Dialog

The HappyAssistant algorithm capitalizes on a major advantage of natural language, which is the ability to ask very general questions and elicit rich descriptive responses from the user (i.e., the ability to choose from a very large set of possible attributes). Obtaining a larger set of user requirements through the dialog shortens the interaction length and thus improves its quality. The dialog is initiated by matching concepts from the user's query to business rules. If a match is found, then a web page associated with that rule is presented to the user. Otherwise, based on the weights of partially matched rules, the system finds the most important missing concept and asks a question to elicit descriptive responses from the user.

More specifically, each rule has a rank. Rank is defined as the number of concepts in a rule that has *not* been extracted from user queries in a particular session. For example, if a rule requires three concepts to trigger, and only two of those concepts have been identified, then this particular rule is assigned a rank of one. The initial matching algorithm iterates through all the rules in the knowledge base, calculates the rank of each rule. The presence of any rules with rank zero indicates the triggering of that rule, i.e., that particular product or service matches the user description. In this case, a detailed web page regarding that item is displayed in a separate browser window. If there are multiple rule triggerings, the rule with the highest weight is selected. If there are no rules of rank zero, then additional refinement is needed to recommend a suitable item. HappyAssistant chooses from the rules of the lowest rank and the highest weight, and from that finds a concept that has not yet been identified. The system retrieves the natural language question associated with that concept from the Domain Lexicon and poses it to the user, prompting for a response. Simultaneously, all items associated with partially matched rules (up to a certain adjustable upper limit) are offered to the user as 'items of interest.' The user has the option to either answer the question posed or browse through the list of relevant items.

The natural language question posed by the HappyAssistant for each concept is not merely intended to affirm or deny that concept, but to further elicit descriptive responses from the user. This allows the system to collect more information about user needs. For example, the question associated with the affordable-concept (people who are looking for an affordable computer) might be "What are your financial constraints?" This tactic would allow the HappyAssistant to differentiate among similar concepts such as value-performance-combination (i.e., a mixture of value and performance) or

unconstrained-finances (i.e., don't care about price) in a single exchange. After the concepts are extracted from the user reply, the partially matched rules will be re-matched again. This refinement process of question and answer repeats until the system can recommend an item to the user or until all possible items have been eliminated. (In which case, a graceful exit message is displayed.)

To assist users who have trouble describing their requirements, the HappyAssistant implements a back-off mechanism that rephrases the question in a more specific way. For example, if the question "What are your financial constraints" does not elicit any recognizable response, the system would proceed to ask a different question, "Is affordable price important to you?"

In general, this dialog based information access task is accomplished by implementing a hybrid forward and backward chaining rule based system. An initial description is gathered from the user, creating a set of currently active rules, which is refined by subsequent query and response. Because the rules operate on concepts, they can be independently developed apart from the language analysis section of the system.

2.4 Other Natural Language Processing Components

The language analysis in the prototype system is limited. We currently only apply techniques in noun phrase parsing, and simple language generation. The components in the dotted boxes in Figure 1 are placeholders and will be addressed in the next version of the system.

The noun phrase parser is used to process grammatical and semantic information of interest from the user input. This shallow parser extracts the head of the noun phrase from its modifiers and identifies key concepts of interest that are present in the user's query. For example, if the user is looking for a "thinkpad with external mouse", the parser identifies the "laptop" and "mouse" concepts and marks the "laptop" as the headword, which represents the item the user is interested in. On the other hand, if the user enters "mouse for laptops," the parser identifies the same concepts, but marks "mouse" as the headword. Parsing, as opposed to pure keyword matching, improves performance by distinguishing between head and attributes and identifying key concepts rather than keywords.

The natural language generation is applied when the system recommends products to the user. In the prototype system, when a product is recommended, an explanation of the recommendation will be automatically generated. It integrates the concepts detected in the user's utterance from the history of interactions.

2.5 Screen Shots and a Walkthrough

Several screen shots are attached in the Appendix to show an example of interactions. Screen 1 provides a text field for the user's NL query. A cue "I'm looking for" is given in the anticipation that users will follow with descriptive noun phrases. In this example, the user types in the query "a notebook for my consulting business". The PM applies the Noun Phrase Parser and detects the "thinkpad" and "business-use" concepts. The DM receives two concepts and checks it with the Knowledge Base for business rules. Some rules are partially matched. The DM sends these potential rules to the AM. The AM shows a list of interesting items corresponding to the potential rules, together with brief descriptions for the user to browse (as in Screen 2). In addition, the DM compares ranks and weights of the potential rules and prompts a question "Please describe your final constraints" for more information. The user can browse the interesting items as he/she wishes or can choose to answer the question to narrow down the search space. In Screen 2, the user types in "not important, but the performance is essential". Based on this input, the PM discards "affordability" concept and detects "high-performance" concept. The DM matches the previously active rules and sends a further narrowed set of potential rules to the AM. The AM will retrieve a list of high performance models and show to the user. Once again, by comparing ranks and weights of the rules, the DM prompts another question "are you looking for something that is top of the

line?” (as in Screen 3) The answer of “yes, absolutely” triggers “best-performance” concept and has a complete match with a business rule for Thinkpad 770 (in Screen 4). Based on the interaction history, a paragraph (summary) is generated to explain to the user why that product is recommended. The user can follow the “here” link or the picture icon for the additional information about Thinkpad 770.

3 Comparative Evaluation

We conducted a user study to evaluate the natural language dialog based prototype system in comparison with a fully developed (by an independent organization) menu driven system. For the user testing, we addressed the following questions: Can natural language based navigation be more efficient (number of clicks, time spent searching, etc.) and easier to use than menu driven navigation? By how much? What are users' responses toward natural language based navigation as opposed to menu driven navigation? How do users with different online experiences react to the natural language dialog based navigation?

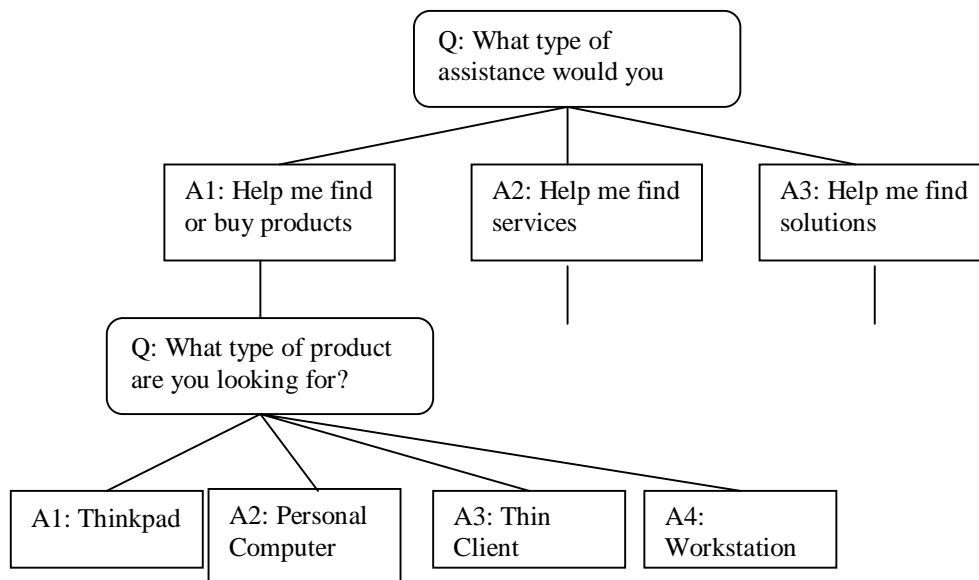


Figure 2: Example of Menu Navigation Structure

3.1 Menu Driven System

Menu driven systems are commonly used navigation devices on business oriented Web sites. Typically, menu-driven systems offer a limited number of options to choose from by displaying radio buttons, choice boxes, or pull-down menus. The system used in the user testing is based on a question answer paradigm. The system prompts to the user a question with a list of answers. When the user chooses one answer, the system will provide another question with a list of answers. The process continues until the system exhausts all predefined questions/answers and reaches the final recommendation. An example of this kind of structure can be found in Figure 2.

3.2 Testing Background

An independent testing agency recruited a total of 17 people. They tested the natural language dialog based system and the menu driven system. A screener was used to recruit the participants. Among those participants, four of them considered themselves with advanced computer skills, eight with intermediate level of proficiency and five with limited experiences with internet.

A testing room was set up with a division to allow for the respondent and the moderator to work at one monitor while we were positioned behind the divider at another monitor. The systems were linked in such a way that we could manipulate the Happy Assistant prototype if necessary and observe the testing as it occurred without interfering with the interview. Such manipulation was intended only for some fatal errors like infinite loops or unexpected exceptions, and was rarely used. In the testing, it turned out less than 5% of interactions were intervened.

Each interview began with an introduction by the moderator explaining the purpose of the interview. It was explained that they would be using two prototype web sites. In addition, they were informed of the moderator’s independent and objective position and encouraged to be open with their opinions of the prototype. After the introduction, the participants were given various scenarios. These scenarios were designed to let them experience critical parts and navigation of each web site in order to form an opinion of the tool’s concept. They were then asked to rank the scenarios on a 1 to 10 scale (where 10 is easy) with regards to the ease of navigation and the series of events leading up to the result. The moderator probed the participant throughout the scenario, and the participants were asked for their overall reaction to the concept of each prototype upon completion. In total, six scenarios were used for the testing. For each scenario, there were two similar versions presented: one for the natural language system and the other for the menu driven system. Each participant was randomly assigned three scenarios.¹

3.3 Observations and Results



Figure 3: Number of Clicks and Average Time Spent to Accomplish Each Task

Comparing NL dialog based navigation with the menu driven navigation in finding products, the number of clicks is significantly reduced by 63.5% (indicated by the T-test, $P < 0.0005$) and the amount of time spent is significantly reduced by 33.3% (indicated by the T-test, $P < 0.025$). The comparison can be found in Figure 3. The horizontal axis is the scenario number and the vertical axis represents the metrics of interest. Reactions from users with different online experiences varied. The less experienced users preferred the NL enabled navigation much more than the experienced users. After each task, we asked participants to rate the ease of use of the two systems. The average rating of subjects with limited internet

¹ An example of a scenario is: one version: “Your daughter’s birthday is coming soon. Currently she is a high school student and you would like to get her a computer for birthday. You would also like to have internet access at home and occasionally play some computer games with your kids,” and the other version: “You have just moved into a new house and would like to get a computer for your home. You are interested in internet access to check email, browse websites, and trade online from home. Occasionally, you would also need to write letters.”

experience was 9.4 for the natural language based system and 6.3 for the menu driven system. The average rating of users with intermediate experience was 8.5 for the natural language system and 8.1 for the menu driven navigation. The average rating of users with advanced experience was 8.3 for the natural language system and 8.9 for the menu driven navigation. Figure 4 shows the average ratings for each scenario by three different groups. Since the scenarios were randomly assigned, it turned out that no one with limited experience tested scenario 5 and no one with advanced experience tested scenarios 1 and 6.

In the testing, respondents preferred the NL dialog based navigation (HappyAssistant) to the menu driven navigation two to one (2:1). Respondents thought the HappyAssistant was extremely easy to use, and they were comfortable and confident with the resulting information it provided. Users liked the fact that they could express their needs in their “own jargon” instead of the foreign “computer jargon”. There was

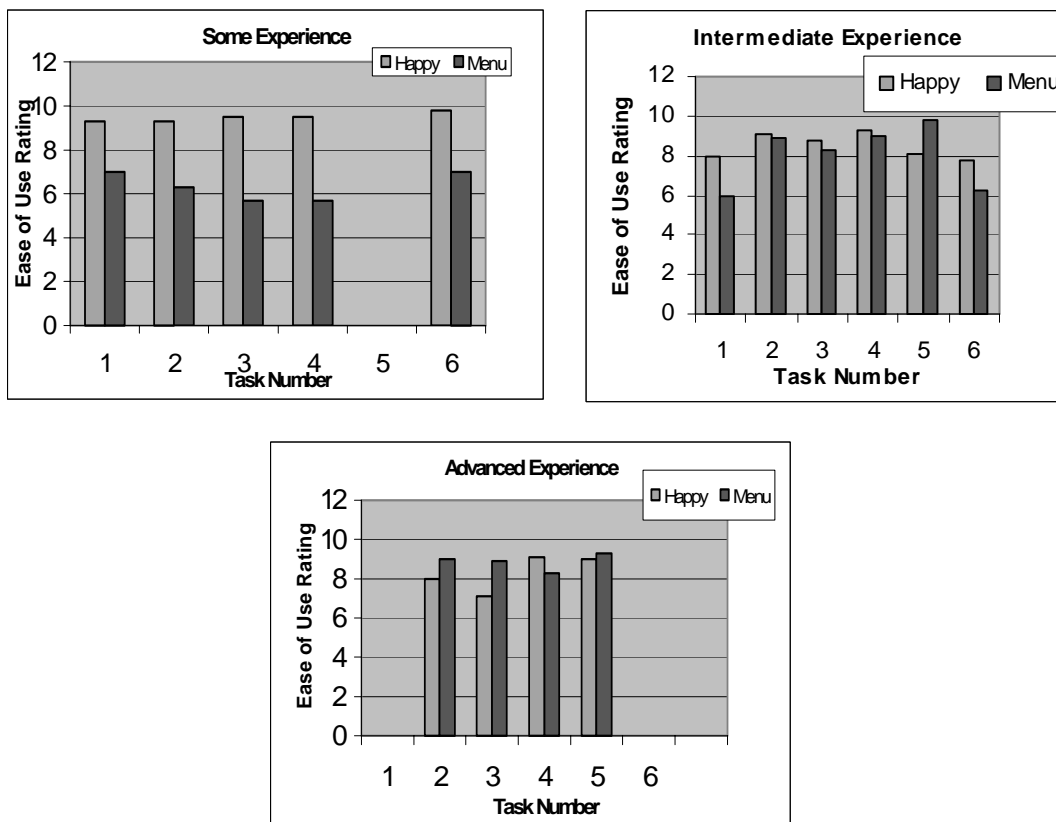


Figure 4: Ratings from Users with Different Levels of Experience

also the perception that with the HappyAssistant model, the computer did all the work for them instead of them doing all the work for the computer (as in the menu-driven model).

Respondents were extremely pleased with the personal search instead of the generalization of the individual’s needs. Many commented positively on the pronouns used by the site such as “*I have found something for you*”. The personalized care makes them “feel like an individual” and not like “70,000 other people”. Moreover, many respondents preferred to type in what they were looking for and not to answer generic questions that did not apply to their needs.

Despite the fact that most users preferred the NL dialog based navigation, there was great evidence of the utility of menu driven searches. There were definitely users who liked the ability to select options from a menu, specifying that the multiple-choice method was easy. There were also users who liked having questions asked of them. Typically, such users were either not comfortable with their typing ability or unable to express what they are looking for without additional information.

3.4 Discussion

In this study we have found that natural language dialog can shorten the time spent on the interaction, and provide a more natural interface that users prefer, especially less experienced ones. Furthermore, we found the sophistication in dialog management is more important than the ability to handle complex natural language sentences. Although the moderator explicitly explained to the participants that they could express their interests in natural language, 85% of input was in the form of keywords or noun phrases. We have learned that the current internet keyword search engines have created a “search culture” which is widely accepted by most internet users. Some users felt it is quicker to type in just phrases or key words; and others were concerned with typing and spelling. This observation suggests that complex sentence analysis might not be a key issue in the e-commerce environment.

We have also learnt that in order to improve the functionality of an e-business site, the natural language dialog navigation and the menu driven navigation should be combined to meet users’ different needs. While the menu driven can provide choices for the user to browse around or learn some information, the natural language dialog provides the efficiency, flexibility and natural touch to the users’ online experience.

Moreover, in designing NL dialog based navigation, one important issue is to show users that the system does understand his/her requests before giving any recommendation or relevant information. This can be achieved by summarizing the user’s requests by paraphrasing it using context history, or by engaging in meaningful conversations with the user. The study showed that almost all of the respondents appreciated the additional questions prompted by their input and the summary coming with each recommendation.

4 Conclusions and Future Work

This paper describes a prototype system that provides natural language dialog capabilities to help users access e-commerce sites to find relevant information about products and services. The prototype system leverages technologies in natural language processing and human computer interaction to create a faster and more intuitive way of interacting with websites, especially for the less experienced users. The result of a comparative study shows that users prefer the natural language enabled navigation two to one over the menu driven navigation. In addition, the study confirmed the efficiency of using natural language dialog in terms of the number of clicks and the amount of time required to obtain the relevant information. Comparing to the menu driven system, the average number of clicks used in the natural language system was reduced by 63.2% and the average time was reduced by 33.3%.

We have learned that users like natural language dialog mode of interaction. In the context of e-commerce, the natural language dialog is particularly useful for understanding user’s intentions and providing additional information and recommendation. Furthermore, to provide easy access to information on e-commerce sites, natural language dialog based navigation and menu driven navigation should be intelligently combined to satisfy user’s different needs.

The work presented in this paper is a proof of concept study. Although the prototype system has employed only basic techniques in natural language processing and human computer interaction, the results learnt from the user testing are significant. By comparing this simple prototype system with a fully

deployed menu system, we have learned that users, especially novice users strongly prefer the natural language dialog based system. We have also learned that in an e-commerce environment, sophistication in dialog management is more important than the ability to handle complex natural language sentences. In the prototype system, we have explored the issue of mapping common sense concepts to the business specifications through business rules. We believe that, in order to successfully facilitate the natural language dialog, the translation between common sense concepts to business ontologism is very important. Our current and future work, extending the results of this paper, includes enhancing the natural language analysis and dialog management modules and automatically learning business rules and the ontological mapping between customer terms and business terms.

References:

- Aggarwal, C., Wolf, J., and Yu, P. (1998), A framework for the Optimizing of WWW Advertising, *Trends in Distributed Systems for Electronic Commerce, LNCS 1402*, Lamersdorf and Merz Eds.
- Androutopoulos, I. and Ritchie, G. D. (1995), Natural Language Interfaces to Databases – an Introduction, in *Natural Language Engineering 1(1)*: 29-81, Cambridge University Press.
- Bray, T., Paoli, J. and Sperberg-McQueen, C. M. (1998), Extensible Markup Language (XML) 1.0., Technical Report, <http://www.w3.org/TR/REC-xml-19980210>, World Wide Web Consortium Recommendation.
- Carpenter, B. and Chu-Carroll, J. (1998), Natural Language Call Routing: A Robust, Self-organizing Approach, in *Proceedings of the Fifth International Conference on Spoken Language Processing*.
- Chu-Carroll, J. and Carpenter, B. (1998), Dialog Management in Vector-based Call Routing, in *Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics*.
- Huberman, B. A., Pirolli, P. L. T., Pitkow, J. E. and Lukose, R. M. (1998), Strong Regularities in World Wide Web Surfing, in *Science*, Vol. 280.
- Muller, J. and Pischel, M. (1999), Doing Business in the Information Marketplace, in *Proceedings of the 1999 International Conference on Autonomous Agents*, Seattle, USA.
- Radev, D., Kambhatla, N., Ye, Y., Wolf, C. and Zadrozny, W. (1999), DSML: A Proposal for XML Standards for Messaging Between Components of a Natural Language Dialog System, in *Proceedings of the AISB'99 (Artificial Intelligence and Simulation of Behavior) Workshop on Reference Architecture and Data Standards for NLP*, Edinburgh, England.
- Saito, M. and Ohmura, K. (1998), A Cognitive Model for Searching for Ill-defined Targets on the Web – The Relationship between Search Strategies and User Satisfaction, in *Proceedings of 21st International Conference on Research and Development in Information Retrieval*, Australia.
- Walker, M., Fromer, J., and Narayanan, S. (1998), Learning Optimal Dialogue Strategies: A Case Study of a Spoken Dialogue Agent for Email, in *Proceedings of 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics*, Montreal, Canada.
- Zadrozny, W., Wolf, C., Kambhatla, N. and Ye, Y. (1998), Conversation Machines for Transaction Processing, in *Proceedings of the Fifteenth National Conference on Artificial Intelligence (AAAI) and Tenth Conference on Innovative Applications of Artificial Intelligence Conference (IAAI)*, Madison, Wisconsin, USA.

Appendix



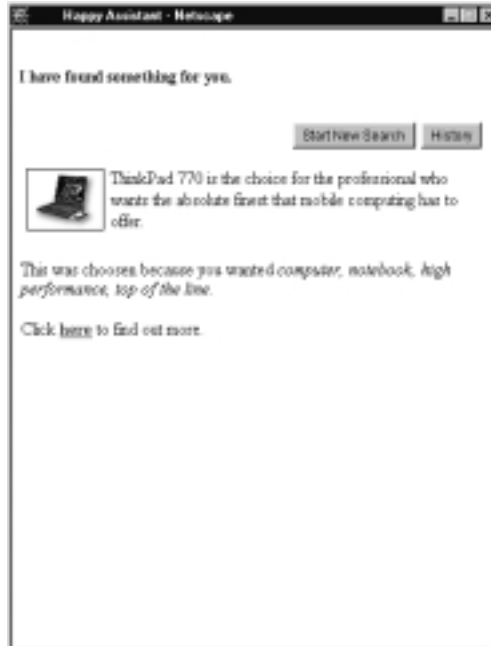
Screen 1



Screen 2



Screen 3



Screen 4

