

# Domain Metric Knowledge Model for Embodied Conversation Agents

Goh Ong Sing, Chun Che Fung, Arnold Depickere and Kok Wai Wong  
School of Information Technology, Division of Arts,  
Murdoch University, Perth, Western Australia 6150, AUSTRALIA  
{os.goh | l.fung | a.depickere | k.wong} @ murdoch.edu.au

*Abstract*— In this paper, we propose a domain metric knowledge model for Embodied Conversational Agents (ECAs) for human computer interaction. In our research, we have found that by incorporating Open-Domain and Domain-Specific knowledge, a new generation of ECAs will be able to demonstrate intelligence and to increase trust by the users. The proposed system is aimed at advancing the scholarly field of artificial intelligence and to offer assistance to queries on real-world issues. The experimental system we have developed aims to provide answers to questions on pandemic threat. The focus is on Bird Flu as the Domain-Specific knowledge together with the support of multiple open-domain knowledge bases. We also present a scheme for selecting specific websites as the input knowledge bases for our conversational system based on reputability, credibility, reliability and trustworthiness.

*Index Terms*— Embodied Conversational Agents (ECAs), Knowledge Model, Artificial Intelligence (AI), Information Extraction (IE)

## I. INTRODUCTION

With the rapid growth of data and information in many digital libraries and databases, there is a need for a new generation of computational theories and tools to assist human in extracting useful information or knowledge from these repositories. To have an appreciation of the size of these information storages and as an example, Google<sup>1</sup> has more than 14 billion pages. While it seems that they are ready sources of knowledge to be used by conversational systems, the sheer volume of information calls for changes on how to utilize the World Wide Web online documents as a source of knowledge bases. This source of information is usually generated dynamically and automatically after a site is updated. On the other hand, Web metadata provides the topology of a site, which includes information on the neighboring pages, leaf nodes and entry points. In addition, the availability of the Internet has made it possible to build

knowledge bases by mass collaboration, with thousands of volunteers contributing to the knowledge bases simultaneously.

In this paper we propose a Domain Metric Knowledge Model for Embodied Conversation Agents (ECAs). True intelligent action requires large quantities of knowledge. The inability to acquire appropriate and sufficient knowledge has long been the major constraint preventing the rapid adoption of conversational systems. Using manual approaches to input handcrafted knowledge is time consuming, costly and impractical. In our proposed approach, we aim to harvest the vast reservoir of knowledge from the internet by deploying the domain metric knowledge bases architecture. This will assist the constructions of large-scale knowledge bases to be used as the engine for the future intelligent conversation systems.

## II. EMBODIED CONVERSATIONAL AGENTS

An Artificial Intelligent Neural-network Identity (AINI) architecture has been developed and reported in [1]. AINI can be scaled up and applied to new knowledge domain and to handle queries concerning specific topic of interest such as Bird Flu pandemic. AINI engine is portable and has the ability to disseminate information to the public through conversational dialogues. It is also capable to carry on multiple independent conversations at the same time. AINI's knowledge bases use plug-in principle which can quickly be augmented with specific knowledge to help targeted groups.

The reported project involves the establishment of a Crisis Communication Network (CCNet) portal<sup>2</sup>. The objective is to use the embodied conversational agent (ECA) to provide meaningful and quick responses to queries concerning the pandemic in a conversation. The completed system is aimed to advance the knowledge and applications of ECA. The research also contributes to the discipline of artificial intelligence and to offer practical assistance to those who seek answers from this site. Our real-time prototype is based on distributed agent architecture designed specifically for the Web and human-computer communication

---

<sup>1</sup> Google hits on 11/10/2006 was 14,210,000,000 by given “the” in the input query.

---

<sup>2</sup> CCNet portal can be access at <http://ainibot.murdoch.edu.au/ccnet>

systems. The software system consists of conversation engine, multidomain knowledge model, multimodal human-computer communication interface and multilevel natural language query. The communication protocol is based on TCP/IP for its universal acceptance. It is believed that AINI as a conversation agent is fulfilling its mission in maintaining a meaningful conversation with users who interact with AINI.

In this paper, we present a domain metric knowledge model for our conversation agent. With the amalgamation of advanced computer technologies such as computer animation, computer generated graphics and artificial intelligence, we are getting closer to the possibility of creating realistic virtual personalities which are capable to converse with human in natural language. Our proposed architecture involves a talking virtual embodied character that is capable of holding a natural conversation with the user. While the user types in the input, our agent responds using a speech on demand system use Text-to-Speech (TTS) technology with lip sync encoder. During the past decade, rapid advances in spoken language technology, natural language processing, dialogue modeling, multimodal interfaces, animated character design, and mobile applications all have stimulated interest in a new class of conversational interfaces [2], [3], [4], [5]. In contrast with command or dictation-style interfaces, conversational interfaces aim to support spontaneous conversation with a

large vocabulary repository. The ability to utilize large amount of vocabulary is essential in maintaining an ongoing dialogue between a user and computer. They also permit the computer to be able to respond appropriately to the user's topic or queries. In order to accomplish these goals, next-generation conversational interfaces requires supports by parallel natural language processing for both input (e.g., speech recognition) and output (e.g., TTS). Currently, our system is not only capable to communicate through the Web. This system is also possible to be implemented on the mobile networked devices such as mobile phone, Palm and Pocket PC devices[6].

### III. DOMAIN METRIC KNOWLEDGE MODEL

A significant difference between our proposed research agent and other conversational agents is our Domain Metric Knowledge Model as shown in Figure 1. In our approach, we define the knowledge base of our conversation system is a collection of specific conversation domain units. Each unit handles a specific knowledge used during the conversation between the user and the computer. For example, in the Open-Domain knowledge, the subject domains will cover subjects such as personality, business, biology, computer, etc. In this report, our focus is on the medical subject and in particular pandemic bird flu.

Domain knowledge model plays a major role in

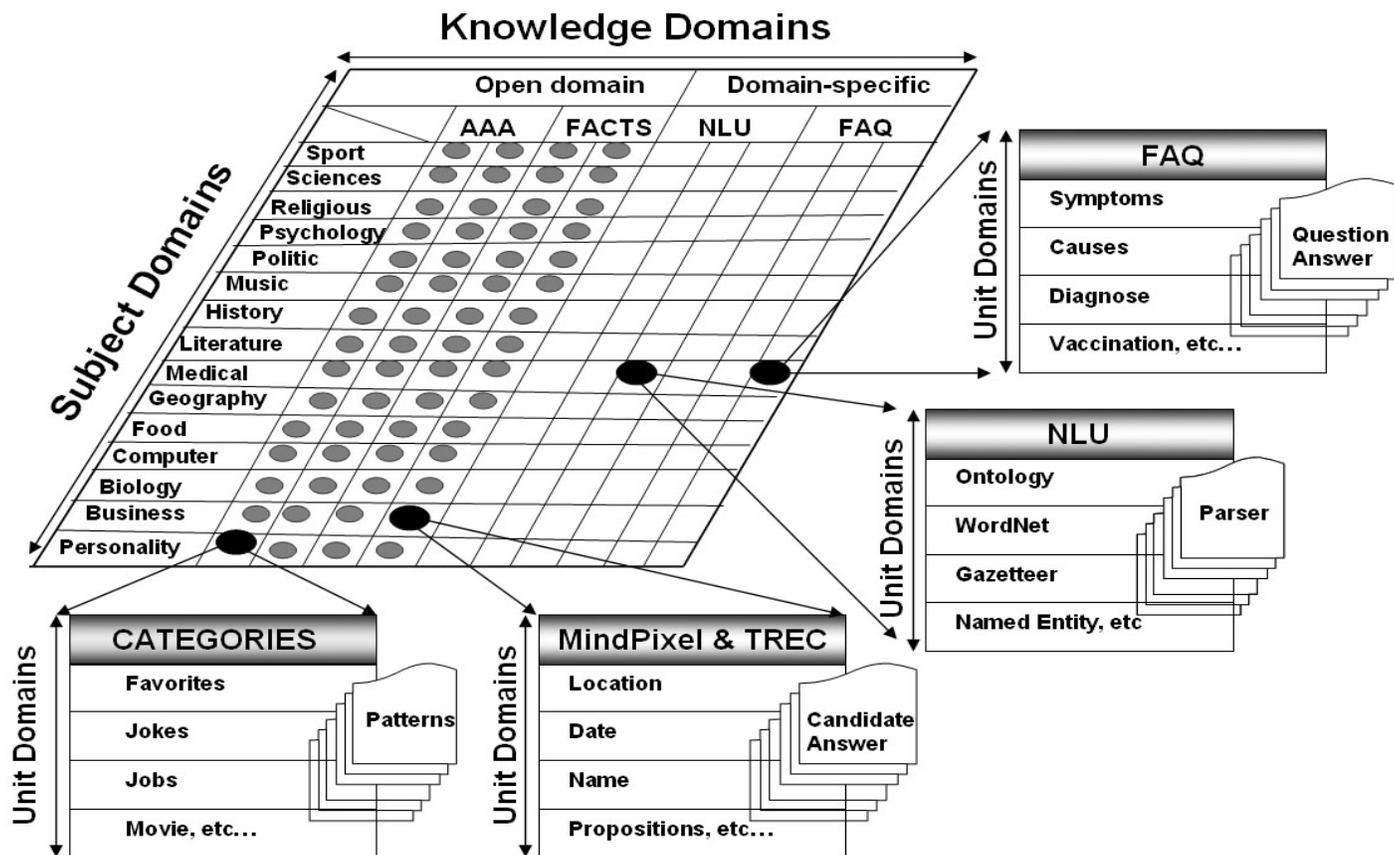


Fig. 1 Architecture of the Domain Metric Knowledge Model for ECAs

conversational agent systems. Systems that rely on specially crafted knowledge-base are normally composed of two subcategories: the *traditional* or *narrow domain*, and the *open domain*. In the traditional domain, systems attempt conversational fluency based on limited domains of expertise. ELIZA for example simulated a Rogerian psychotherapist, and this implementation is commonly known as DOCTOR. The Rogerian psychotherapist knowledge base attempts to engage the patient by repeating the patient's statements and to encourage the patient to continue talking [7]. Terry Winograd's SHRDLU[8] is another program simulates a robot which is able to interact within a simple world which consists of coloured building blocks and boxes on a flat surface. SHRDLU has knowledge about his world and it can answer questions about it in natural language. However, it is anticipated that SHRDLU will not be competitive enough to win the Turing test because of its limited domain of discourse.

Three decades have passed since ELIZA was created. Computers have become significantly more powerful while storage space and memory size have increased exponentially. These advances have given researchers the opportunity to provide embedded human-like knowledge bases for ECAs. We follow the assumption that human's brain consists of world knowledge, but it has limitations on memory retrieval. Hence, knowledge needs to be specified. In our system, the domain metric knowledge model is designed along this line. Our system uses custom domain-oriented knowledge bases and existing knowledge bases from online documents and training corpus. It contains a spectrum of possible solutions to queries on specific domains to general conversation questions. In our study, the domain model is the taxonomy of knowledge related to the specific topic. It can also be considered as XML-like metadata model. This approach will reduce the need to predict every possible input from the user. Instead, this approach allows the manager of the system to put more effort on how to handle conversations within a specified domain or the domain-specific conversations.

In this project, the most important contribution is the development of the "*domain knowledge plug-in components*". By this arrangement, the specific domain knowledge becomes portable, scaleable and easy to be incorporated in other applications. In short, this domain model can be considered as domain-dependent modular components. This approach will allow future improvements to encourage openness and collaborative contribution to the specific knowledge domain. By domain, we refer to the theme that determines the focus of the conversational system. An Open-Domain will enable the development of a wide range of information sources. For a system that focuses on certain domains, it is more likely that the techniques are more restrictive and logic-based. There will be relatively limited available sources as compared to an Open-Domain system. In the Open-Domain system, candidate answers are ranked according to individual features such as how well the answer

matches the question. A domain-oriented conversational system deals with questions within a domain-specific environment will be seen as a richer approach. This is because natural language processing systems can exploit domain knowledge and ontologies. Advanced reasoning such as providing explanations for answers, generalizing questions, etc however is not possible in Open-Domain systems.

#### *A. Open-Domain Knowledge Bases*

Open-Domain conversational systems need to deal with questions about nearly any topic. It is very difficult to rely on ontological information due to the absence of wide and yet detailed world knowledge. On the other hand, these systems have much more information and data to be used in the process of answering the queries. In AINI's conversation system, we deployed the large-scale mass collaboration Mindpixel [9] and training data sets from Text REtrieval Conference (TREC) training corpus [10]. It also uses ALICE Annotated AIML (AAA) [11], the award winner Loebner Prize[12] and the Chatterbox Challenge Winner [13].

Mindpixel is a common sense component and it is similar to OpenMind<sup>3</sup> and Cyc<sup>4</sup>. The system accepts public contributions. However, Cyc model and OpenMind had a bottleneck which prevent truly large-scale collaboration [14]. MindPixel started collecting their propositions privately via email on 1994 and then it evolves to online mass collaboration. To date the project's user base of nearly fifty thousands people has contributed more than one million propositions and recorded almost ten million individual propositional response measurements. AINI's uses only 10% of the Mindpixel propositions. In practice, 10% of the training corpus is held back from training to act as a generalization test to ensure the system did not simply memorize the corpus. Passing this generalization test would be the basis to claim that the system is able to replicate human-level intelligence in a machine. Although a lot of knowledge has been collected, it is recognized that the system is still less than the hundred-of-millions to billions of "pieces of knowledge" that are estimated to be involved with human intelligence [15].

Second common sense knowledge deployed by AINI is a training corpus from TREC. TREC, organized each year by NIST, has offered since 1999 a specific track to evaluate large-scale open-domain QA systems. Finding textual answers to open-domain questions in large text collections is a difficult problem. In our system, we only extracted factoid questions to be incorporated in the AINI's engine. Our concern is the types of questions that could be answered in more than one way. We have tried to avoid such questions. In conversational systems, factoid questions should have only one single factual answer. This will be considered as a good stimulus-response type of knowledge unit. Examples of such questions are, "*Who is the author of the book,*

---

<sup>3</sup> [www.openmind.org](http://www.openmind.org)

<sup>4</sup> [www.cyc.com](http://www.cyc.com)

*The Iron Lady: A Biography's of Margaret Thatcher?*", *"What was the name of the first Russian astronaut to do a Spacewalk?"* or *"When was the telegraph invented?"* TREC corpus has considerably less answer redundancy than the web and thus, it is easier to answer a question by simply extracting the answers from the matching text. To gather this data, we automatically classified questions in the TREC 8 through TREC 10 test sets by their wh-word and then manually distinguished factoid question to representing around half of the initial corpus as shown in the Table 1.

Table 1: Factoid Question from TREC 8, 9 and 11

TREC	Factoid Question	Text Research Collection
8	196	<ul style="list-style-type: none"> <li>Financial Times Limited (1991, 1992, 1993, 1994)</li> <li>the Congressional Record of the 103<sup>rd</sup> Congress (1993), and the Federal Register (1994)</li> <li>Foreign Broadcast Information Service (1996) and the Los Angeles Times (1989, 1990).</li> </ul>
9	692	Set of newspaper/newswire documents includes: <ul style="list-style-type: none"> <li>AP newswire</li> <li>Wall Street Journal</li> <li>San Jose Mercury News</li> <li>Financial Times</li> <li>Los Angeles Times</li> <li>Foreign Broadcast Information Service</li> </ul>
11	109	MSNSearch logs donated by Microsoft and AskJeeves logs donated by Ask Jeeves.

The third knowledge base in the AINI's Open-Domain knowledge model is obtained from The Annotated ALICE AIML (AAA)<sup>5</sup>, a award winner Loebner Prize[12] ALICE chatbot knowledge base. AAA is a free open-source software based on XML specifications. It is a set of Artificial Intelligence Markup Language (AIML) scripts comprising the award winning chat robot. The AAA is specifically reorganized to make it easier for conversational system developer to clone the ECA's brain and to create custom conversation agent personalities, without having to invest huge efforts in editing the original AAA content. AAA knowledge bases covered a wide range of subject domains based on the bot's personality. Example subjects include AI,

games, emotion, economics, film, book, fiction, sport, stories, science, epistemology and metaphysics.

ALICE won the 2000, 2001 and 2004 Loebner Prize for being the most lifelike machines. The competition uses the Turing test, named after British mathematician Alan Turing, to determine if the responses from a computer can convince a human into thinking that the computer is a real person. In the competition, ALICE used a library of over 30,000 stimulus-response pairs written in AIML. The development of ALICE is based on the fact that the distribution of sentences in a conversations tend to follow Zipf's Law[16]. It is indicated that the number of "first word" is only limited to about two thousand. The frequency of any word is roughly inversely proportional to its rank in the frequency table. The most frequently used word will occur approximately twice as often as the second most frequent word. It in turn occurs twice as often as the third most frequent word, and so on so forth. Starting with "WHAT IS", tend to have Zipf-like distributions. This type of analysis which used to require Dr. Zipf many hours of labor is now accomplished in a few milliseconds of computer time. While the possibilities of what can be said are infinite, what is actually said in conversation is surprisingly small. Specifically, 1800 words cover 95% of all the first words input. It is this principle that AINI is operating on.

#### B. Domain-Specific Knowledge Bases

At present, the World-Wide Web provides a distributed hypermedia interface to a vast amount of information available online. For instant, Google [17] currently has a training corpus of more then one trillion words (1,024,908,267,229) from public web pages. This is valuable for our type of research. The Web is a potentially unlimited source of knowledge data; however, commercial search engines are not the best way to gather answers from queries due to the overwhelming of results from a search.

Before the rise of domain-oriented conversational agents based on natural language understanding and reasoning, evaluation is never a problem as information retrieval-based metrics are readily available for use. However, when ECAs began to become more domain specific, evaluation becomes a real issue [18]. This is especially true when Natural Language Understanding (NLU) is required to cater for a wider variety of questions and at the same time to achieve high quality responses.

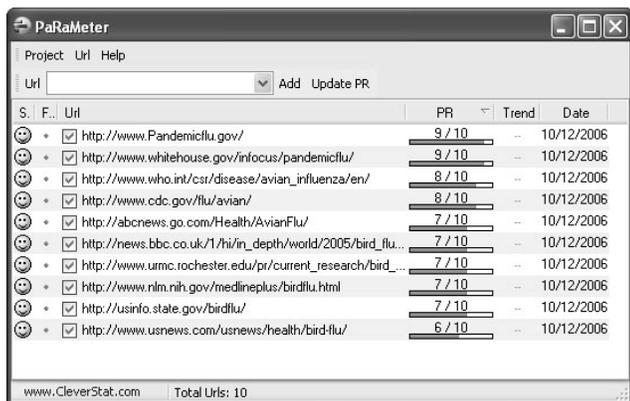
As shows in the Figure 1, AINI's Domain-Specific knowledge base consists of Natural Language Understanding (NLU) and Fequently Asked Question (FAQ). Both components are extracted from the online documents using an Automated Knowledge Extraction Agent (AKEA)[19]. It was aimed at providing up-to-date information to its users via AINI. Another objective of AINI is to deliver essential information from trusted sources and it should be capable to interact with its users. The idea is to rely on a human-like communication approach thereby providing a sense of familiarity and ease.

As the Web that we know today becomes increasingly chaotic, over-powering and untrustworthy, selection of trusted

<sup>5</sup> [www.alicebot.org/aiml/aaa/](http://www.alicebot.org/aiml/aaa/)

Web pages is becoming an important factor contributing to its long-term survival as a useful global information repository. This is to avoid rumors, hoaxes and misinformation. In our experiment, the selection of the website are based on PageRank™[20]. algorithm. The core software technology and algorithm, PageRank™, is a system for ranking web pages developed by Larry Page and Sergey Brin at Stanford University. PageRank™ relies on the uniquely democratic nature of the web by using its vast link structure as an

for the second run. Finally, we retrieved 1,428 URLs out of 1,500 URLs related to the domain being investigated. The reduction in number is due to the duplicated and broken link URLs are removed. Based on the 1,428 URLs, we sent a query to Google's PageRank™ directory using PaRaMeter Tool [22] to determine their rankings. Figure 2 shows the results of the top 10 site based on the PageRank™ scale. The PageRank™ scale goes from 1 to 10. A less important site is one with a PageRank (PR) of 1. The most referenced and supposedly important sites are those with a PR of 7 or 10.



S.	F.	Url	PR	Trend	Date
1	✓	http://www.Pandemicflu.gov/	9 / 10	--	10/12/2006
2	✓	http://www.whitehouse.gov/infocus/pandemicflu/	9 / 10	--	10/12/2006
3	✓	http://www.who.int/csi/disease/avian_influenza/en/	8 / 10	--	10/12/2006
4	✓	http://www.cdc.gov/flu/avian/	8 / 10	--	10/12/2006
5	✓	http://abcnews.go.com/Health/AvianFlu/	7 / 10	--	10/12/2006
6	✓	http://news.bbc.co.uk/1/hi/in_depth/world/2005/bird_flu...	7 / 10	--	10/12/2006
7	✓	http://www.umc.rochester.edu/pr/current_research/bird_...	7 / 10	--	10/12/2006
8	✓	http://www.nih.nih.gov/medlineplus/birdflu.html	7 / 10	--	10/12/2006
9	✓	http://usinfo.state.gov/birdflu/	7 / 10	--	10/12/2006
10	✓	http://www.usnews.com/usnews/health/bird-flu/	6 / 10	--	10/12/2006

The final set of URLs was further culled down to include only selected sites with a regulated authority (such as a governmental or educational institution) that controls the contents of the sites. Once the seed set is determined, each URLs page is further examined and rated as either reliable or reputable. This selection is reviewed, rated and tested their connectivity with the trusted seed pages. From this exercise, *whitehouse.gov*, *pandemicflu.org*, *cdc.gov* and *who.int* were selected due to their PageRank™ scale scores which are more than 7. The most important factors in determining the “reliable authority” of a site is based on its history and the number of backlinks to the governmental and international organization links. The more established and relevantly linked will be considered as “stronger” or “more reliable”. This effectively gives the linked site “trust” and “credential”. The selected URLs are then used as the source knowledge base for AKEA to extract the contents on bird flu so as to build AINI’s domain-specific knowledge base.

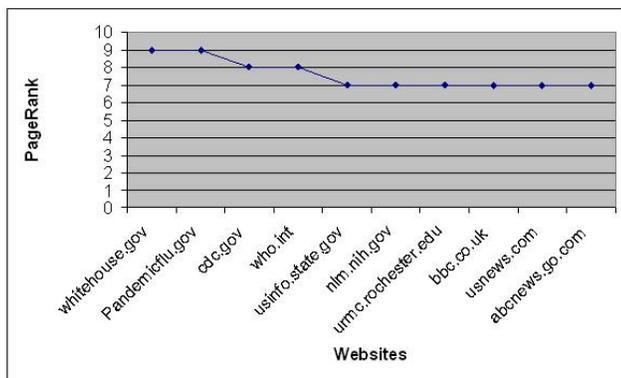


Figure 2. Top 10 Pagerank™ Scale for Bird Flu Domain

indicator of an individual page's value and high-quality sites always receive a higher PageRank™.

In our study, the selection of websites starts with the initial 6 seed words: bird, flu, avian, influenza, pandemic and H5N1. These seeds are supposed to be representative of the domain under investigation. The seed terms are randomly combined and each combination is used in Google API<sup>6</sup> and BootCat Tool [21] for bootstrapping corpora and terms from the web. We use the seeds to perform a set of Google searches, asking Google to return a maximum of 20 URLs per query and get the first corpus. After visual inspection of the corpus, we used top 40 seeds from extract token frequencies

In AINI’s domain knowledge metric, the unit domains in the NLU component consists of knowledge and information harvested or expressed in ontology, gazetteer, named entities and WordNet. These have been implemented as domain-dependent modular components. The named entity module identifies named locations, named persons, named organization, dates, times and key measures in text form. The information are obtained by AKEA. We have also developed a new system for diseases based on symptoms, causes, diagnosis, vaccination locations, persons and organizations. In order to identify these entities, our system uses rules in which we specify the named entities structure in term of text tokens and what can be found about them from resources such as tagger, morphosyntactic analyzer and knowledge bases of names, clue words and abbreviations.

The web knowledge base is then continuously updated with facts extracted from online *pandemic news* using information extraction (IE) by AKEA. IE is the task of extracting relevant fragments of text from larger documents and to allow the fragments to be processed further in automated ways. Application example is this is to prepare an answer for a user’s query. The ontology and gazetteer has been implemented as domain-dependent modular components which will allow future improvements and to maintain openness in the domain knowledge.

In terms of AINI’s FAQ component, the unit domain consists of information concerning diseases, symptoms, causes, diagnosis, vaccinations, etc. The selected FAQ trusted Web page have been

<sup>6</sup> <http://www.google.com/apis>

carried out using PageRank™ as discussed above. But in this stage, each of the selected websites was evaluated in order to find the most suitable and reliable FAQ pages. From this experiment, *pandemicflu.gov* and *who.int* websites have been selected as the source of information for AKEA to build AINI's FAQ knowledge base.

Based on the proposed approach, the quality of the results returned from AINI's engine using the FAQ knowledge base are either similar or better to those generated by search engines such as Google. AINI SQL engine uses the most significant words as keywords or phrase. It attempts to find the longest pattern to match without using any linguistic tools or NLP analysis. In this component, AINI does not need a linguistic knowledge unit and relies on just a SQL query. All questions and answers can be extracted from the complete database which was built by AKEA after applying a filtering process to remove unnecessary tags. Results from our work are discussed in the next section.

## VI. RESULTS AND CONCLUSION

In our study, the experiment on the Domain Metric Knowledge Model has shown interesting results from the natural conversation agent. The key assumption is that important queries do not necessarily turn up the answers that they can be found in a single domain but they may come from different domains. As shown in Table 2<sup>7</sup>, currently AINI's Open-Domain knowledge has more than 150,000 entries in the commonsense stimulus-response categories. Out of these, 100K came from Mindpixel, 997 factoid questions from the TREC training corpus and 45,318 categories from AAA knowledge bases. On the domain-specific, AINI's had more than 1,000 online documents extracted using AKEA. This makes up 10,000 stimulus response categories in total. AINI also has 158 FAQ pairs of question-answers which have been updated using AKEA. In addition, AINI also has collected more than 382,623 utterance conversations with online users since 2005. These utterances will be updated to the AINI's knowledge bases through supervised learning by domain expert. At present, AINI has learnt 50,000 categories from conversations with online users. All of these combined knowledge has made up the total number of the 206,473 stimulus response categories in AINI's Knowledge bases. The original and simple conversational programs such as ELIZA[23] written by Professor Joseph Weizenbaum of MIT has only 200 stimulus response categories. ALICE Silver Edition was ranked the "most human" computer has about 120,000 categories which include 80,000 Mindpixel.

<sup>7</sup> The stimulus response categories of the AINI's knowledge bases calculated on 11 October 2006 and can be access at <http://ainibot.murdoch.edu.au>

Table 2: AINI's Knowledge Bases

Knowledge Bases	Sources	Categories
Domain-Specific	NLU	10,000
	FAQ	158
Open-Domain	MindPixel	100,000
	TREC Corpus	997
	AAA	45,318
Supervised Learning	Conversation	50,000
	Logs	
<b>TOTAL</b>		<b>206,473</b>

As illustrated by the discussion in this paper, there are still many technical challenges to be overcome for conversational characters to be able to scale-up and to become viable in real-world applications. There are a number of open questions that we will continue to explore. Most conversational engines have been developed with the goal of finding topically relevant documents from variety of domain knowledge. Finding accurate answers will be the ultimate goal for the conversational system and it will require different matching infrastructure. We are currently exploring on how best to generate snippets for use in answer mining. Finally, timing is another interesting and important issue. We noted how a correct answer to some queries may change over time. For example, it would take a while for a news or Web collection to correctly answer a question like "*What is the next disease outbreak after bird flu?*" We are making good progress and we shall continue to explore other issues in this interesting discipline of research.

## REFERENCES

- [1] Ong Sing Goh, "A Global Internet-based Crisis Communication: A Case study on SARS using Intelligent Agent," Kolej Universiti Teknikal Kebangsaan Malaysia, Malacca, Malaysia, Technical Report PJP/2003/FTMK(1), January 2005 2005.
- [2] S.L Oviatt, C. Darves, and R Coulston, "Toward Adaptive Conversational Interfaces: Modeling Speech Convergence with Animated Personas," ACM Transactions on Computer-Human Interaction (TOCHI), vol. 11, 2004.
- [3] Victor Zue and James Glass., "Conversational interfaces: Advances and challenges," Proceedings of the IEEE, vol. 88, pp. 1166--1180, 2000.
- [4] J. Cassell, T. Bickmore, M. Billinghurst, L. Campbell, K. Chang, and H Vilhjalmsjon, "Embodiment in conversational interfaces: Rea," presented at the CHI 99 Conference on Human Factors in Computing Systems, New York, 1999.
- [5] Oliver Lemon and Alexander Gruenstein, "Multithreaded context for robust conversational interfaces: Context-sensitive speech recognition and interpretation of corrective fragments," ACM Transactions on Computer-Human Interaction (TOCHI), vol. 11, pp. 241-267, 2004.
- [6] Ong Sing Goh, C. C Fung, Cemal Ardil, K.W Wong, and A Depickere, "A Crisis Communication Network Based on Embodied Conversational Agents System with Mobile Services," Journal of Information Technology, vol. 3, pp. 257-266, 2006.

- [7] Joseph Weizenbaum, Computer Power and Human Reason: W.H. Freeman and Company, 1976.
- [8] Terry Winograd, Understanding Natural Language: Academic Press, 1972.
- [9] mindpixel, (2006),GAC-80K Mindpixel, [Online]. Available: <http://www.mindpixel.com/chris/gac80k-06-july-2005.html>
- [10] NIST, (2006),Text retrieval Conference (TREC), [Online]. Available: <http://trec.nist.gov/>
- [11] Alicebot, (2006),The Annotated A.L.I.C.E. AIML, [Online]. Available: <http://www.alicebot.org/aiml/aaa/>
- [12] Hugh Loebner, (2006),Loebner Prize, [Online]. Available: <http://www.loebner.net/Prizef/loebner-prize.html>
- [13] Chatterboxchallenge, (2006),ALICE Winner of Chatterbox Challenge 2004, [Online]. Available: <http://www.chatterboxchallenge.com/>
- [14] Matthew Richardson and Pedro Domingos, "Building large Knowledge Bases by Mass Collaboration," presented at K-CAP'03, Sanibel Island, Florida, USA, 2003.
- [15] Erik Mueller, (2001),Common Sense in Humans, [Online]. Available: <http://www.signiforum.com/erik/pubs/cshumans.htm>
- [16] Wentian Li, (2006),Zipf's Law, [Online]. Available: <http://www.nslj-genetics.org/wli/zipf/>
- [17] Alex Franz and Thorsten Brants, (2006),All our N-gram are Belong to You, [Online]. Available: <http://googleresearch.blogspot.com/2006/08/all-our-n-gram-are-belong-to-you.html>
- [18] Ong Sing Goh, Cemal Ardil, Wilson Wong, and C. C Fung, "A Black-box Approach for Response Quality Evaluation Conversational Agent System," International Journal of Computational Intelligence, vol. 3, pp. 195-203, 2006.
- [19] Ong Sing Goh and Chun Che Fung, "Automated Knowledge Extraction from Internet for a Crisis Communication Portal," in First International Conference on Natural Computation. Changsha, China: Lecture Notes in Computer Science (LNCS), 2005, pp. 1226-1235.
- [20] Monica Bianchini, Marco Gori, and Franco Scarselli, "Inside pagerank. ACM Transactions on Internet Technology," ACM Transactions on Internet Technology., vol. 5, 2005.
- [21] M Baroni and S Bernardini, "bootcat: Bootstrapping corpora and terms from the web," presented at Fourth Language Resources and Evaluation Conference, 2004.
- [22] cleverstat, (2006),parameter, [Online]. Available: <http://www.cleverstat.com>
- [23] J. Weizenbaum, "ELIZA - A computer program for the study of natural language communication between man and machine," Communications of the ACM, vol. 9, pp. 36-45, 1966