

# Voice-Based Humanoid Robot Interaction

El Sayed M. Saad, Medhat H. Awadalla, Hosam Eldin I. Ali, Rasha F. A. Mostafa

**Abstract**—recently, the interest in service robots endowed with communicative capabilities has been increased. These robots should operate in cluttered and uncluttered environments and interact with humans using natural language to perform a variety of service-oriented tasks. Recognizing and fetching of a user-specified object can be considered as one of the major tasks for a humanoid robot. To get the robot capable of identifying the geometric shapes and colors of the objects, a vision system is proposed. Furthermore, the paper proposes a natural language understanding system, where the robot will be able to effectively communicate with human through a dialogue developed in Arabic language. The developed dialogue and a dynamic object model are used for learning the semantic categories and object descriptions. In this paper, a robot will be developed to interact with the users performing some specified actions. Moreover, integration between the proposed vision and natural language understanding systems has been presented. Finally, a voice-based dialogue between the user and robot will be developed. Intensive experiments have been conducted indoor to address the validity of the complete proposed system. The achieved results show that the overall system performance is high compared with the related literature to the theme of this paper.

**Index Terms**— Vision System, Speech system, object category recognition, Object Detection, Color detection, Natural Language Understanding, Ontology, Syntax, knowledge Representation, Semantic Networks.

## I. INTRODUCTION

Humans are the most advanced creatures of the nature. It is believed that humanoid robots will be the most advanced creatures of humans. Among the man-made creatures such as automobile hand-phones and multimedia devices, robots of future will hopefully be the most ideal assistants to human beings [1]. In the future we will see "personal robots" that will entertain, comfort and serve people in their private lives and homes. While presently robotic servants or butlers exist only in the form of early prototypes in a few research laboratories, they are expected to become as ubiquitous as PCs in the future [2-4]. The necessary functions for the task are understanding users' commands, recognizing user-specified objects by vision, and manipulating objects [5].

An important aspect of humanoid robots in a natural environment is the ability to acquire new knowledge through learning mechanisms, which enhances an artificial system with the ability to adapt to a change or new environment. In contrast to the most offline learning algorithms applied in machine learning today, online algorithms need to be performed automatically, and through interaction with the

environment or with other agents/humans. Here, in this paper, the proposed vision system and dialogue offer the appropriate means. The fact that robots have to be autonomous in such a way that they should do everything without the intervention of humans. Since the proper system is the good vision system, so the question arises here is: how to develop a robot that can see like a human? For many applications in robot vision interested in locating the object by giving it a distinctive color from the surrounding environment as an application to recognize the ball in pitch between two teams of humanoid soccer robot team [6, 7], using laser, sonar, or using camera for robot vision system, or learning to classify objects into categories in human development. Such ability is crucial for robots that should operate in human environments where object categorization skills are required to recognize complex object categories (e.g., metal objects, empty bottles, etc.) [8]. However in this paper, robots will learn how to distinguish among different geometric shapes of square, rectangular, circle, and triangular objects picked up via a camera mounted on the robot and also identify their colors.

The paper also focuses on the distinction of robot to a command given by user in Arabic language. Using the Arabic language syntax for imperative sentence and establishment of dialogue to identify the objects that does not exist in the database. In this paper, we address learning of unknown objects in dialogue, which enables a robot to acquire information about unknown objects, and stores this information in a knowledge base. A typical problem will be raised is that non-trivial information must be communicated, such as when the user enters an imperative syntax error, or there are new words in the written sentence that cannot be understood by the system. Thus, the dialogue system should conduct dialogue strategies for learning in such way that the information about the object can successfully be communicated. In addition, it has to cope with new words, grammatical and semantic levels to achieve the learning goal. It should create a model of the object's semantics, which describes the type, color, shape, properties of the object and its function. All previous data will be addressed using Arabic language. Moreover, this paper adds question syntax of Arabic language. Finally, using a special tool is developed to convert any answer of robot to Arabic speech. The remainder of this paper is organized as follows: Section 2 gives an overview of the proposed system comprising the proposed vision system architecture, the natural language understanding system architecture, and the integration of them. Section 3 presents experiments and discussion. Section 4 concludes the paper.

## II. THE PROPOSED SYSTEM

The proposed system, as shown in figure (1), consists of the following sub-systems:

1. Vision Sub-system.

**Manuscript received September 02, 2012.**

**El Sayed M. Saad**, Department of Communications, Electronics and Computers, Helwan University, faculty of Engineering Cairo, Egypt.

**Medhat H. Awadalla**, Department of Communications, Electronics and Computers, Helwan University, faculty of Engineering Cairo, Egypt & Department of Electrical and computer engineering, SQU, Oman.

**Hossam E.I.Ali**, Department of Communications, Electronics and Computers, Helwan University, faculty of Engineering Cairo, Egypt.

**Rasha F.A. Mostafa** Department of Communications, Electronics and Computers, Helwan University, faculty of Engineering Cairo, Egypt.

2. Natural Language Understanding Sub-system.
3. Merging Sub-system.
4. Talking Sub-System.

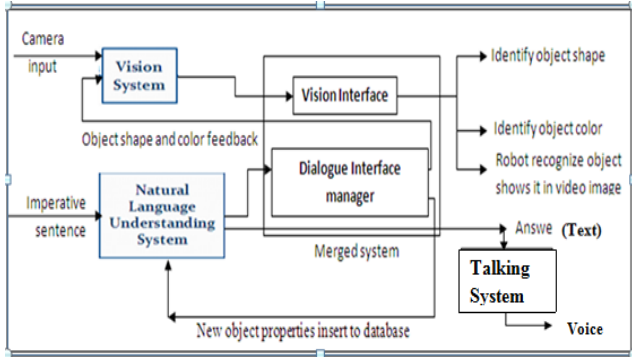


Fig.1. System overview

**1. Vision Sub-system**

There are things that attract any child such as colors and geometric shapes of objects, once the child taught how to distinguish between them using one word that defines each shape and color. He could recognize them by himself later on. The main aim of the proposed system is to make robots behave as child, once it is learned the skill of how to recognize object’s shape such as square, rectangle, circle, or triangle, and its color such as red, black, white, blue, green,... etc., seen by its camera. It could detect the object’s shape and its color by itself seen later in any image taken by his camera at different places. The proposed vision system has the following procedure and is demonstrated in the flowchart shown in figure2.

The vision system procedure:

**A- Acquiring image**

Images taken from the robot’s camera can easily be fed to Matlab program using the ‘*videoinput*’ function. This function makes it possible to assign a variable as a video input. Image processing cannot be performed on a video input, so single frames should be extracted from the video with a frame grabber [9]. A so-called snapshot is taken out of the video input and this single image is used for object and color detection.

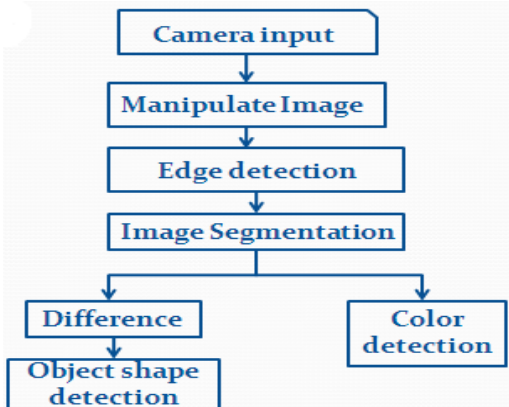


Fig.2. The object and color detection program flow chart

**B- Image manipulation**

Image frame taken from camera is a color image. In Matlab, images automatically are coded using RGB-space. In the RGB color space, each color is described as a combination of three main colors, namely Red, Green, and Blue. This color space can be visualized as a 3d matrix. Each image is

converted into black and white then the image is filtered to remove any added noise due to lighting. A well-known noise filter is the median filter. In Matlab, this filter can be used with the ‘*medfilt2*’-function [9], as shown in figure 3.



Fig.3. result of converting image to black and white

**C- Edge detection**

The black and white image is converted into edge image, as shown in figure (4), using ‘*edge*’-function [9].



Fig.4. result of edge detection

**D- Image Segmentation**

The previous image is used to find the boundary of each object by using ‘*boundaries*’-function [9]. These indices are used to cut the black and white image to a set of images, each of which contains only one item, as shown in figure 5, each image is then used to determine the object’s shape and color.

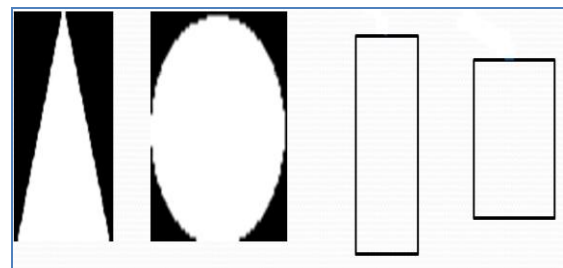


Fig.5. The segmented object image

**E- Identification of the shape of an object**

There are many ways to determine the location of the element, including the mean and variance of number of ones in an image. However, the previous methods cannot determine the shape of the object; accordingly, a new method is proposed that calculates the difference between the number of ones in each row. If it always increases, the object’s shape will be a triangle, or if it increases before the middle line and then decreases after the middle line of the object, the object’s shape will be a circle, and if it does not change, the object is a square or a rectangle, it depends on the dimensions of the image as shown in figure 6. All written functions are based on ‘*diff*’-function [9].

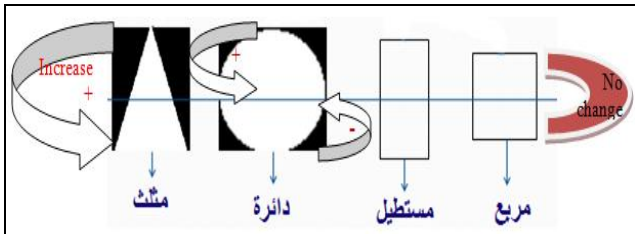


Fig.6. difference of each object

**F- Identification of the color of an object**

Identification of the color of the object can be achieved through the borders of each object, which it is obtained previously and separated into an individual image, and then the color can be determined by creating a color map using 'colormap' - function [9], where each color can be visualized as a 3d matrix. Finding the average of this matrix can identify the color of the object from the table (1) shown below.

Table. I Colormap

	B	G	R
Black	0	0	0
White	1	1	1
Red	0	0	1
Blue	1	0	0
Green	0	1	0
Cyan	1	1	0
Magnetic	0	1	1

The previous procedure has been applied to the image shown in fig. 7, and the achieved results confirm the validity of the proposed approach. All results identified the color and shape of the object and it is given to the user in syntax of Arabic Language.

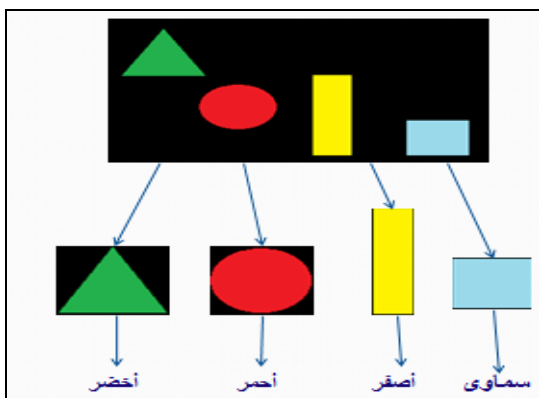


Fig.7.color identification

**2. Natural Language Understanding System Sub-system**

An interactive learning for artificial systems has been addressed in several systems. However, the number of approaches that allow interactive knowledge acquisition for humanoid robots is still comparably small [10]. This paper focuses on how to establish a dialogue between the user and the robot especially if some of the commands to the robots are not pre-defined. Furthermore, the paper concentrates on how the robot will understand the commands on syntax of the Arabic language, to address these issues, the natural language understanding system is proposed as shown in figure 8. The proposed system has the following components:

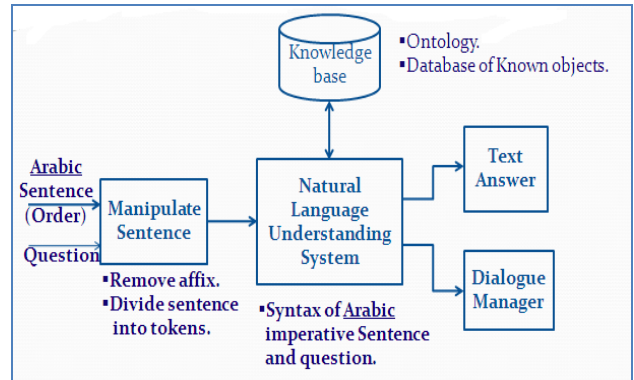


Fig.8. Natural Language Understanding System Overview

**A- The form of imperative sentence syntax**

First sentence is entered into the system in the form of imperative syntax; imperative syntax of the Arabic language takes more than one form as shown in figure9.a. The input sentence is divided into a set of tokens, and then the affix such as "ال" added at the beginning of the word or "ى" added at the end of the word is omitted. Moreover, the proposed system has ability to take a question from user about the place of the object using a question tool of "أين". The question is entered into the system in the form of question syntax; the question syntax of the Arabic language takes more than one form as in the shown in figure 9.b, it begins with "أين".

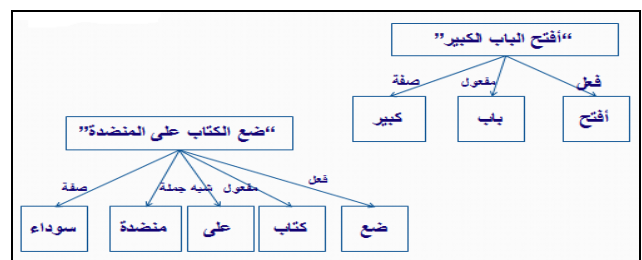


Fig.9.a. Example of the form of imperative sentence syntax

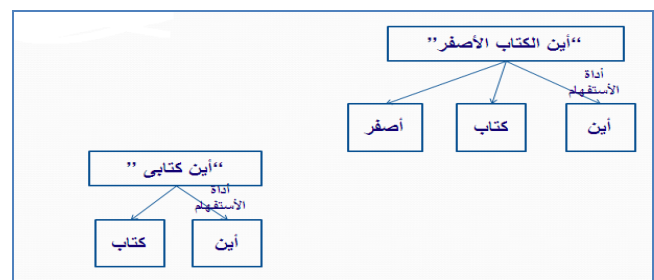


Fig.9.b. Example of the form of question syntax

**B- Ontology**

Our ontology inspired from the work in [10, 11], however it is applied to the Arabic language, as shown in figure 10, and also we have added that the object is classified by its color and shape. Knowledge representation is an area in artificial intelligence that focuses on the design of formalisms which can explicitly represent knowledge about a particular domain, and the development of reasoning methods for inferring implicit knowledge from the represented explicit knowledge.



Semantic network form a family of knowledge representation formalisms which can be used to represent and reason with conceptual knowledge about a domain of interest. For the classification of each object and storing its data in the database, we used semantic network[12].A semantic network is a simple representation scheme that uses a graph of labeled nodes and directed arcs to encode knowledge [13,14].Information type and semantic categories of objects are modeled in ontology. The object ontology provides inheritance information and defines properties that can be associated with objects. Moreover, our ontology for objects' locations inspired by the work in [15],however it is applied to Arabic language as shown in figure 11.User should store objects' locations found in the environment in the database using "بجانب", "على", "فى", "تحت", "خلف", "or", "أمام" as shown in figure 10.b.

C- Dialogue

Dialogue begins when the robot does not find the object in its database, or the user enters an imperative syntax error. Dialogue is in a form of questions and answers words, mutual between the user and robot in Arabic as shown in figure 11 and it is different from the work in [10]where the answer is yes ,or no only, and in English. There are ten commands the user can use them such as, "أضطر", "أفتح", "أغلق", "أملئ", and,"أنظر", or dialogue begins when user ask for object's place using question of form "أين.....؟",and the robot does not find the object in its database, the dialogue takes a form as shown in figure 10.a

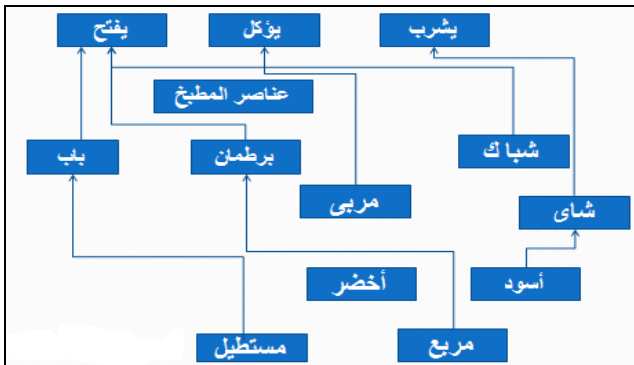


Fig.10.a. Ontology organization with functional concepts, type hierarchy and properties

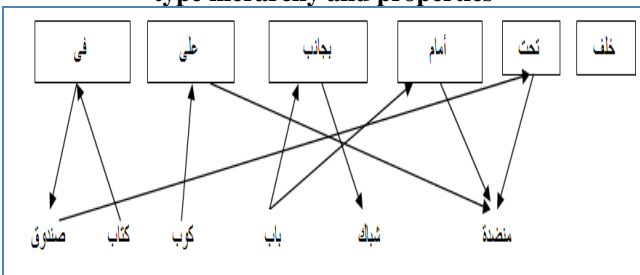


Fig.10.b. Example of the ontology organization for objects locations in database

ما هو أسم الشيء؟ كتاب  
 ما الصفة المميزة له؟ أخضر  
 أين يوجد؟ فوق المنضدة

Fig.11. Dialogue sample

3. Merging Sub-system

Merging between the vision and the natural language understanding systems is implemented by regular usage of GUI using Visual BASIC 6.0, where it appears to the user interface screen. One of the major problems faced most of previous researchers is the linking between Matlab environment and Visual Basic to have a reasonable interface for the users. In this paper, all programs have been developed in Prolog, and the Dynamic Link Library is implemented in such way that it can deal directly with Visual Basic programs. The developed interface gives the user different capabilities to choose as follows:

1- The user can ask the robot to describe the scene in the front of itself in terms of the shapes of the objects and their colors by using the developed vision program just by pressing a key called the vision system in the developed interface shown in figure12.a.

2-The user can ask the robot to identify or fetch a particular object, the robot will use the developed Natural Language Understanding System to search about the object in its database. If it is recognized, the robot writes to the user that object is found. Then the object features will be sent to the vision system to recognize its shape and Color, as shown in figures 12.b and 12.c.

3-The user can ask the robot about the object's place, the robot will use the developed Natural Language Understanding System to search about the object in its database. If it is recognized, the robot writes to the user that the object place as shown in figures 12.e and 12.f. Then the object will be sent to the vision system to recognize its shape and Color.

In all the above cases, if the robot could not recognize the object because it is not in its database or the user entered a syntax error, or, robot does not find object's place in its database, a dialogue between the user and the robot will start, the user will answer some questions as shown in figure 12.d, and then the database will be adapted to accommodate the new information. If the same command sent to the robot, or the robot is asked again for identifying the same object, the system will be able to identify it.

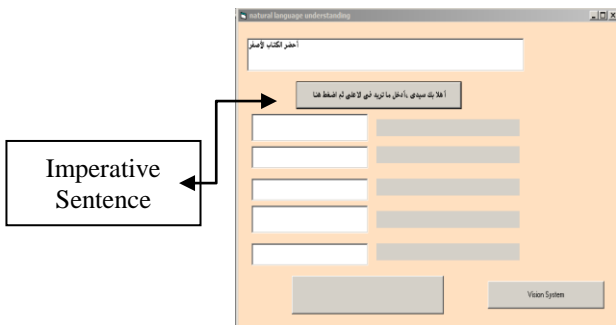


(a) Result from calling vision system

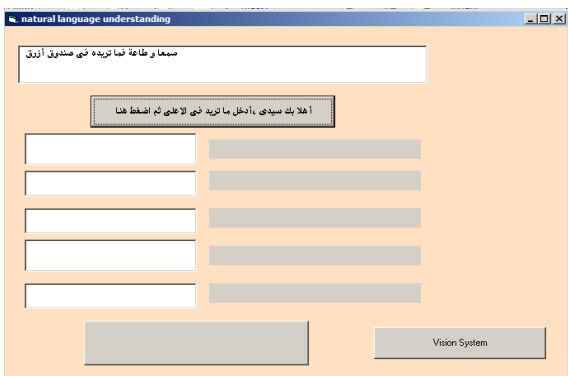
(f) the robot writes to the user the object place  
Fig.12. Example of the proposed system scenario

#### 4. Voice-based interactive dialogue (Talking sub-System)

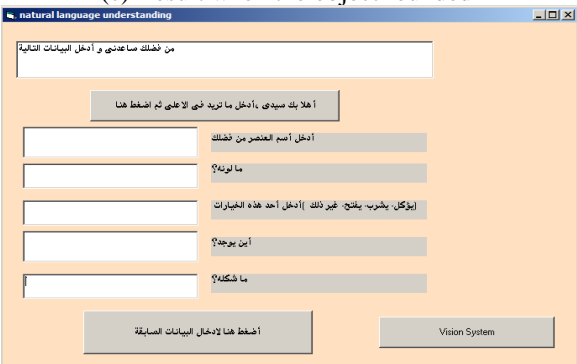
The previously mentioned subsystems are interdependent. It is not enough to equip the robot with basic functionalities for dialogue comprehension and production to make it interact naturally in situated dialogues. We also need to find meaningful ways to relate language to be spoken, and enable the robot to use its perceptual experience to continuously learn and adapt itself to the environment. And cover speech recognition, where the dialogue between the user and robot will be voice based. For this purpose we used a tool to convert Arabic text of robot answer to voice, this tool called MbrolaTools35. This tool is well described in [16]. We used it to convert the answer of robot appeared to user to speech heard by user. When answer is written to user it will hear at the same time, this by pass this text to MbrolaTools35 as seen in Figure (13.a, 13.b), however the screen appears in figure 13.b will be invisible for the user. User will hear the voice immediately when it appears at text box. Also the dialogue between user and robot shown in figure 12.d changed to be not only based on text but also voice based.



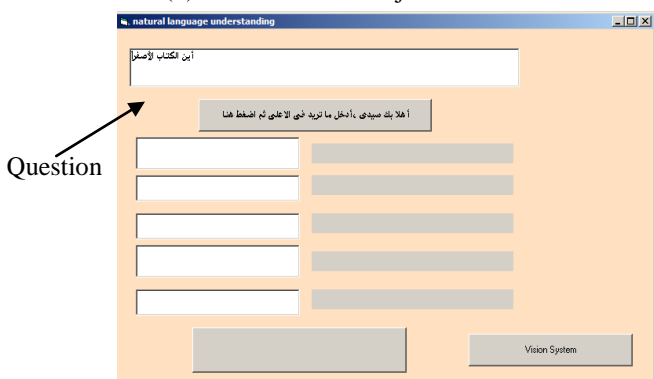
(b) Example of imperative sentence



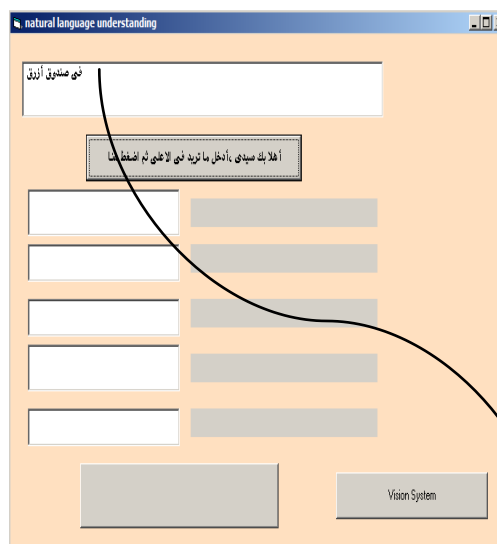
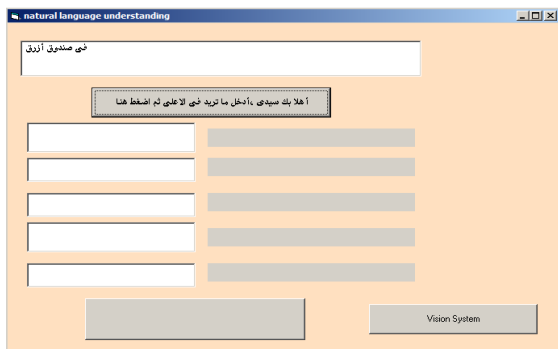
(c) Result when the object founded



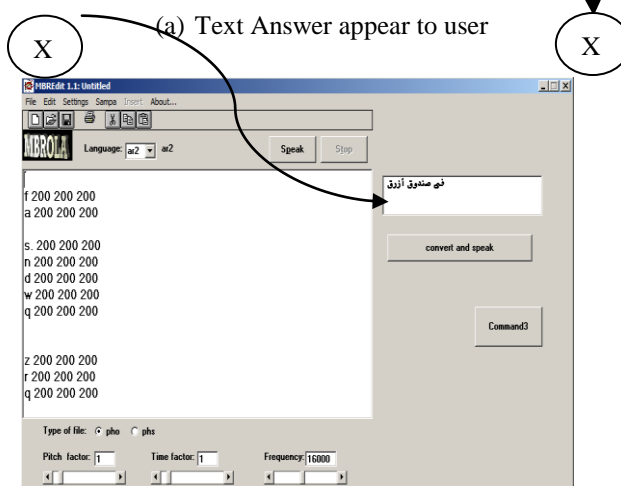
(d) result when the object not founded



(e) The user ask the robot about the object's place



(a) Text Answer appear to user



(b) Converting robot answer from text to speech

Fig.13. Example of converting text to speech mechanism

III. EXPERIMENTS AND DISCUSSION

Intensive experiments have been conducted to address the validity of the proposed systems. First, we have tested the vision program in several stages to check its accuracy, initially it is tested on still images, and then pictures from a camera held on the laptop, the destination in an interview of an embodiment of the kitchen, and the kitchen components, such as door, window, a cup and also a piece of cheese cooked. Background with one color, black, is chosen. The achieved results from the vision program have a remarkable precision as shown in figure 14, although the vision affected by the camera resolution, and lighting. Our camera resolution used was (640\*480). Furthermore, the achieved results show that the accuracy of natural language understanding program is also reasonable even for different users as long as they know the basics of Arabic. We provide a comparison with H.Holzapfel , D.Neubig, A.Waibel [10] as seen in table II.

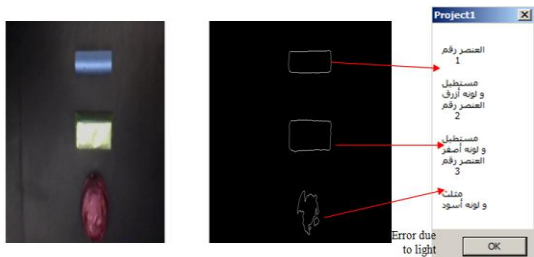


Fig.14 results of a practical image

Table- II Comparison between the proposed system and the work of [10]

	H.holzapfel , d.neubig, a.waibel[10]	Proposed system
<b>Dialogue language</b>	English	Arabic
<b>Dialogue technique</b>	Question and using yes or no answer	Question and using word answer
<b>Ontology</b>	Provide imperative sentence only	Provide imperative sentence and question using where
<b>Vision system</b>	Using software tool Visual object recognition is not the main focus of the paper	Provide new technique to visual object recognition and identify the object features

IV. FUTURE WORK

The work of this paper will be extended so that the robot can determine a suitable place for the user to put things, and also provide a road map to the user towards that place. In addition to, a real robot will be constructed to test the proposed approaches in real environments and applications.

REFERENCES

[1] M. Vukobratović, "Humanoid Robotics, Past, Present State, Future", Director Robotics Center, Mihailo Pupin Institute, 11000 Belgrade, P.O. Box 15, Serbia, E-mail: vuk@robot.imp.bg.ac.yu, SISY 2006 • 4th Serbian-Hungarian Joint Symposium on Intelligent Systems, pp 13-27.

[2] V. Graefe, R. Bischoff, "Past, Present and Future of Intelligent Robots", Intelligent Robots Lab , LRT 6, Bundeswehr University Muenchen, 85577 Neubiberg, Germany, http://www.UniBw-Muenchen.de/campus/LRT6,CIRA 2003, Kobe, pp 1-10.

[3] C.Pasca, "History of Robotics", University of Ottawa, ENRICHMENT MINI-COURSE, Robotics – Intelligent Connection of the Perception to Action, May 5, 2003, pp1-46 .

[4] R. JARVIS, "INTELLIGENT ROBOTICS: PAST, PRESENT AND FUTURE", International Journal of Computer Science and Applications, Vol. 5, No. 3, pp 23 – 35, 2008.

[5] M. Takizawa, Y. Makihara, N. Shimada, J. Miura and Y. Shirai, " A Service Robot with Interactive Vision- Object Recognition Using

Dialog with User - ", Osaka University, 2-1 Yamadaoka, Suita, Osaka 565-0871, Japan, E-mail: shimada@eng.osaka-u.ac.jp, 2003.

[6] H.J.C. Luijten, "Basics of color based computervision implemented in Matlab", Technische Universiteit Eindhoven, Department Mechanical Engineering, Dynamics and Control Technology Group, Eindhoven, June, 2005, pp 1-24.

[7] E. Menegatti, S. Behnke, C. Zhou, " Humanoid soccer robots", Robotics and Autonomous Systems, contents lists available at ScienceDirect, journal homepage: www.elsevier.com/locate/robot, Robotics and Autonomous Systems 57 (2009) 759\_760.

[8] J.Sinapov and Al. Stoytchev, "Object Category Recognition by a Humanoid Robot Using Behavior-Grouped Relational Learning", Developmental Robotics Laboratory, Iowa State University, {jsinapov, alexs}@iastate.edu, 2011, pp 1-7.

[9] Mathworks Matlab Image Processing function list, http://www.mathworks.com/products/image/functionlist.html, 2012.

[10] H.Holzapfel , D.Neubig, A.Waibel, "A dialogue approach to learning object descriptions and semantic categories", Contents lists available at ScienceDirect, Robotics and Autonomous Systems 56 (2008) 1004\_1013.

[11] J. Carbonell, Towards a self-extending parser, in: Annual Meeting of the Association for Computational Linguistics, 1979.

[12] R. Becher, P. Steinhaus, R. Zöllner, R. Dillmann, "Design and implementation of an interactive object modelling system", in: Proceedings of ISR 2006 and Robotik 2006, Düsseldorf, 2006.

[13] M. Khalifa, V. Liu, " KNOWLEDGE ACQUISITION THROUGH COMPUTER MEDIATED DISCUSSIONS: POTENTIAL OF SEMANTIC NETWORK REPRESENTATIONS AND EFFECT OF CONCEPTUAL FACILITATION RESTRICTIVENESS ", Twenty-Sixth International Conference on Information Systems, 2005, pp 221-232.

[14] P. Tanwar , T. V. Prasad, M. S. Aswal, "Comparative Study of Three Declarative Knowledge Representation Techniques", Poonam Tanwar et. al. / (IJCS) International Journal on Computer Science and Engineering Vol. 02, No. 07, 2010, 2274-2281.

[15] S. H'uwel, B. Wrede, and G. Sagerer, "Robust Speech Understanding for Multi-Modal Human-Robot Communication", Faculty of Technology, Applied Computer Science Bielefeld University, 33594 Bielefeld, Germany, 2006.

[16] Al. Ramsay, H. Mansour, " Towards including prosody in a text-to-speech system for modern standard Arabic", Received 13 March 2006; received in revised form 22 June 2007; accepted 22 June 2007 Available online 6 August 2007, Science Direct, Computer Speech and Language 22 (2008) 84–103.

**El Sayed M. Saad** is a Professor of Electronic Circuits, Faculty of Engineering, Univ. of Helwan. He received his B.Sc. degree in Electrical Engineering (Communication section) from Cairo Univ., his Dipl.-Ing. Degree and Dr.-Ing degree from Stuttgart Univ. , West Germany, at 1967, 1977 and 1981 respectively. He became an Associate Prof. and a Professor in 1985, and 1990 respectively. He was an International scientific member of the ECCTD, 1983. He is Author and/or Coauthor of 132 scientific papers. He is a member of the national Radio Science Committee, member of the scientific consultant committee in the Egyptian Eng. Syndicate for Electrical Engineers, till 1 May 1995, Member of the Egyptian Eng. Syndicate, Member of the European Circuit Society (ECS), and Member of the Society of Electrical Engineering (SEE).

**Medhat H. Awadalla** is an associate professor at Communication and Computer Department, Helwan University. He obtained his PhD from university of Cardiff, UK. Msc and Bsc from Helwan university, Egypt. His research interest includes cloud computing, sensor networks, high performance computing and real time systems.

**Hossam Eldin I. Ali** received his B.Sc. degree in Communications & Electronics Engineering, his M.Sc. degree in Computer Engineering, and his Ph.D. degree in Computer Engineering from Helwan University, Cairo, Egypt, in 2000, 2004, and 2009 respectively. He has three published papers at M.Sc. degree, and seven published papers at Ph.D. degree. He is currently a Teacher at Electronics, Communication & Computer Department, Faculty of Engineering, Helwan University.

**Rasha F. A. Mostafa** is an Assistant Teacher at Communication and Computer Department, Helwan University. Received her B.Sc. degree in Computer Engineering, and her M.Sc. degree in computer Engineering from Helwan University, Cairo, Egypt in 2001 and 2008, respectively. Her research interest includes Artificial intelligent, Robot., Security and operating System.

