

# The Design and Implementation of Online Identification of CAPTCHA Based on the Knowledge Base

Yu'e Song, Chengguo Wang, Ling Zhu, Xiaofeng Chen, Qiyu Zhang

► **To cite this version:**

Yu'e Song, Chengguo Wang, Ling Zhu, Xiaofeng Chen, Qiyu Zhang. The Design and Implementation of Online Identification of CAPTCHA Based on the Knowledge Base. 9th International Conference on Computer and Computing Technologies in Agriculture (CCTA), Sep 2015, Beijing, China. pp.92-99, 10.1007/978-3-319-48354-2\_9. hal-01614186

**HAL Id: hal-01614186**

**<https://hal.inria.fr/hal-01614186>**

Submitted on 10 Oct 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



# The Design and Implementation of online Identification of Captcha Based on the Knowledge Base

Yu'e Song<sup>1,2,a</sup>, Chengguo Wang<sup>1,b</sup>, Ling Zhu<sup>3,c</sup>, Xiaofeng Chen<sup>1,d</sup>, Qiyu Zhang<sup>1,e\*</sup>

<sup>1</sup>*Yantai Academy, China Agricultural University, Yantai, 264670, China;*

<sup>2</sup>*School of electrical and information engineering, Beijing Polytechnic College, Beijing, 100042, China;*

<sup>3</sup>*College of statistics, Shandong Institute of business and Technology, Yantai, 264005, China*

<sup>a</sup>*aeaeae623@163.com*, <sup>b</sup>*wangcg@126.com*, <sup>c</sup>*oklab@qq.com*, <sup>d</sup>*cxfg1979@126.com*, <sup>e</sup>*rcraingo@163.com*

## Abstract

The Completely Automated Public Turing Test to Tell Computers and Humans Apart (CAPTCHA) identification is designed to distinguish between computers and humans and it prevents the web application programs from malicious attacks, so it has been applied widely. However, great challenges must be faced with the development of CAPTCHA identification. In order to improve the safety of the professional system, the CAPTCHA online identification based on the knowledge base, which has high security and bases on semantic questions and the professionalization of professional system, is put forward combining with the recessive CAPTCHA. The specific implementation course of the new online identification method is worked out according to the example of animal identification. The application of the verification code is suitable for people who have the corresponding professional knowledge. Because the computer has great difficulty to answer semantic information questions, which are also professional issues, so the new online identification method based on the verification of knowledge has very high security.

**Keywords:** CAPTCHA, online identification, knowledge base, animal

## 1 Introduction

With the rapid development of internet network, security problem of the web application becomes an extremely important issue for us. The HTTP attack based on the form automatically submission is a common way of network attack. According to the HTTP protocol, the attacker can write program to simulate the method of form submission, and submit the abnormal data to site service automatically and rapidly. This constitutes the basic HTTP attacks. An attacker can repeat logging to break a user's password and this will lead to a leakage of users' privacy information. In order to prevent the attacker using program automatic login, Completely Automated Public Turing test to tell Computers and Humans Apart (CAPTCHA) technology has been widely used<sup>[1]</sup>.

The CAPTCHA is a kind of program algorithm to distinguish between computers and humans, so the procedure must be able to generate and evaluate computer test which human can easily pass but not

---

\*Corresponding author: Qiyu Zhang, Yantai Academy, China Agricultural University, No.2006, Coastal middle Road, Gaoxin District, Yantai City, Shandong Province, China. E-mail: rcraingo@163.com

for computers <sup>[2, 3]</sup>. Because the computer cannot solve CAPTCHA question, the user who answer the question can be considered human <sup>[4]</sup>.

In order to protect the network, CAPTCHA has been applied widely, such as preventing spam ads in the blog post, protecting website registration and the E-mail address, online polls, preventing dictionary attacks, the search engine robots, worms and spam, etc.

Since CAPTCHA has been proposed, different research institutions and scholars have developed a variety of CAPTCHA. CAPTCHA has different ways of classification <sup>[5]</sup>. According to the type of information, CAPTCHA can be divided into text CAPTCHA, image CAPTCHA, graphics CAPTCHA, audio CAPTCHA and video CAPTCHA. According to the way of recognition, CAPTCHA can be divided into dominant CAPTCHA and implicit CAPTCHA. According to the interaction, CAPTCHA can be divided into static CAPTCHA and dynamic CAPTCHA. Along with the development of the CAPTCHA, CAPTCHA recognition technology is also developing and some methods have been put forward, such as the matching shape context <sup>[6]</sup>, template matching <sup>[8]</sup> and neural network identification methods<sup>[7]</sup>. This makes the security of the CAPTCHA has a huge challenge. Dynamic CAPTCHA and recessive CAPTCHA have a good security and is the research direction in the future.

The hidden CAPTCHA <sup>[5]</sup> refers to answering the question of the CAPTCHA expressing according to the semantic of CAPTCHA, for example, CAPTCHA system first randomly generates an expression  $(5+3)*9/4$  and requires the user to answer the expression values; CAPTCHA system picks up a few images from the graphics library and users need to rotate the graphics to the right direction. Though artificial intelligence has a rapid development, the computer has much difficulty to answer semantic information questions, so the hidden CAPTCHA is safe.

In this paper, the CAPTCHA technology is studied deeply. Based on the implicit CAPTCHA and combining with the characteristics of professional system, a new kind of CAPTCHA is proposed based on the knowledge base and the security of the system can be effectively improved using the new kind of CAPTCHA.

## **2 Knowledge representation**

In the knowledge base, knowledge representation methods are logical notation, production representation, frame representation and object-oriented representation, semantic representation and the XML representation and representation of ontology <sup>[9]</sup>, etc. According to the characteristics of the CAPTCHA, we choose production knowledge representation description.

Shortliffe firstly introduced the concept of production in the famous expert system MYCIN. The structure IF (E1 & E2 &... & En) THEN A is called the rule. It means that if the logical expression of E1 & E2 &... & En established, the conclusion A is right. The expression E1 & E2 &... & En is called former part of the rule and is any legal logical expressions. It is the prerequisite for reasoning by using the rule. A is called later part of the rule and is the result of reasoning using the rule. <sup>[10]</sup>. The rule knowledge representation has many advantages, such as simple and clear reasoning, the reasoning machine design and implementation is simple and has a good characteristics in some specific application environment, etc.

### 3 The design of CAPTCHA based on the knowledge base

For some professional systems, CAPTCHA can be structured based on knowledge base. Because users have the corresponding knowledge and can reason the related results according to the precondition. Let us use a simple animal identification as an example to illustrate how to construct CAPTCHA.

We give the following rules about animal identification:

IF the animal has hair THEN the animals are mammals

IF the animal has milk THEN the animals are mammals

IF the animal has feathers THEN the animal is a bird

IF the animal can fly AND lay eggs THEN the animal is a bird

IF the animal eats meat THEN the animal is a carnivorous animal

IF the animal has a canine tooth AND claw AND eyes staring at front THEN the animal is a carnivorous animal.

IF the animal is mammals and has claw THEN the animal is a hoof animal.

#### 3.1 The design of the database and table for knowledge base

According to the rules of reasoning above, we designed the rules table, inferences table and synonym table. Rules table save the atomic conditions of precondition, which are the minimum condition of premise condition. The above animal identification rules are in the rules table as shown in Table 1.

**Table 1.** Animal identification rules

<i>Serial number</i>	<i>Rules</i>
1	Have hair
2	Have milk
3	Have feathers
4	Can fly
5	Can lay eggs
6	Eat the meat
7	Have canine tooth
8	Have claws
9	Eye star at the front
10	Have hoof

The result of reasoning is text messages. There are different representations for the same text messages and the computer can't recognize it very well, therefore automatic word segmentation can be used for the results and CAPTCHA. In this process, the word which not be used can be removed and the keywords will be extracted, then we can match the keyword. For Chinese word segmentation, IK Analyzer 2012 can be used. The IK Analyzer is an open source lightweight Chinese word segmentation toolkit based on Java language. In the 2012 version, we support configuring IKAnalyzer. CFG.XML file to expand proprietary dictionary and stop using dictionary and dictionary format is utf-8 without BOM in Chinese text files<sup>[11]</sup>. Stop using words are not really meaning of function words in both English and Chinese<sup>[10]</sup> and can be ignored because they does not affect the understanding of sentence meaning. The stop using dictionaries are built on the basis of the literature [10] and [11]. In order to assist CAPTCHA judgment, two options are increased which must be contained keywords and must

not contained keywords. Meanwhile, in order to reduce the complexity of the system reasoning, the result is made as easy as possible. Inferences table is shown in Table 2.

**Table 2.** Inference table data

<i>Premise condition</i>	<i>Results</i>	<i>Whether the word is segmented</i>	<i>Must contained keywords</i>	<i>Must not contained keywords</i>
1	Mammals	no	no	no
2	Mammals	no	no	no
3	Birds	no	no	no
4,5	Birds	no	no	no
6	Predators	no	no	no
7,8,9	Predators	no	no	no
1,10	Hoofed animals	no	no	no
2,10	Hoofed animals	no	no	no

Synonym of the word in the results is stored synonym table, including Chinese, English and acronyms.

In the MySQL database we design different table structures, which are shown in Table 3, Table 4 and Table 5.

**Table 3.** Rule table

<i>Field</i>	<i>Data type</i>	<i>Dote</i>
Id	int	Automatic numbering, primary key
Rule	varchar(100)	

**Table 4.** Inferences table

<i>Field</i>	<i>Data type</i>	<i>Note</i>
Id	int	
Condition	varchar(100)	
Result	varchar(100)	Automatic numbering, primary key
Segmentation	char(1)	
Key	varchar(200)	
Antonym	varchar(200)	

**Table 5.** Synonym table

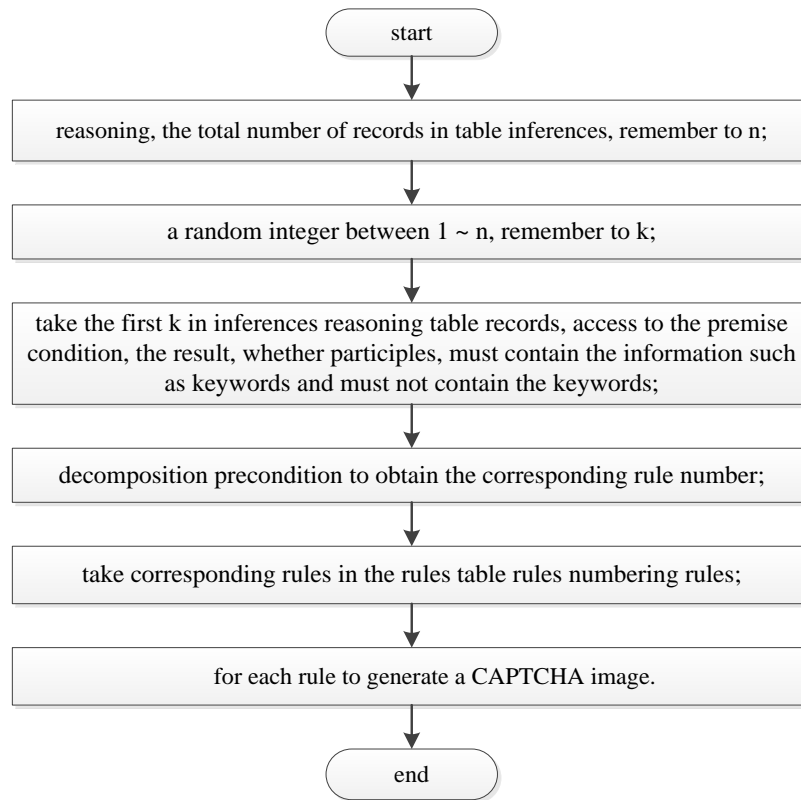
<i>Field</i>	<i>Data type</i>	<i>Note</i>
Id	int	
Key	varchar(100)	Automatic numbering, primary key
Synonym	varchar(100)	

### 3.2 CAPTCHA generation algorithm

- (1) Reason the total number of records in inferences table and remember to n;
- (2) Randomly select the integer between 1 ~ n, remember to k;

- (3) Take the kth records in the inferences table and access the premise condition, the result, whether participles, the keywords which must be contained and which must not be contained;
- (4) Decompose precondition to obtain the corresponding rule number;
- (5) Take corresponding rules in the rules table rules numbering rules;
- (6) Generate a CAPTCHA image for each rule.

The algorithm flow chart is shown in Figure 1.



**Figure 1.** CAPTCHA generation algorithm

### 3.3 CAPTCHA validation algorithm

- (1) The user reasons according to generated CAPTCHA images and enter the CAPTCHA;
- (2) Word segmentation? If no, compare the CAPTCHA entered by the user and the results and judge whether they are consistent. If consistent, agree on. If inconsistent, judge whether there is a synonym and whether consistent after replacement. If unanimity, agree on. If inconsistent, not through;
- (3) If the words need segmentation, do words segmentation to the CAPTCHA entered by the user and results and match the keyword. If they are consistent, agree on. If inconsistent, judge whether there is a synonym and whether consistent after replacement. If unanimity, agree on. If inconsistent, not through. The algorithm flow chart is shown in Figure 2.

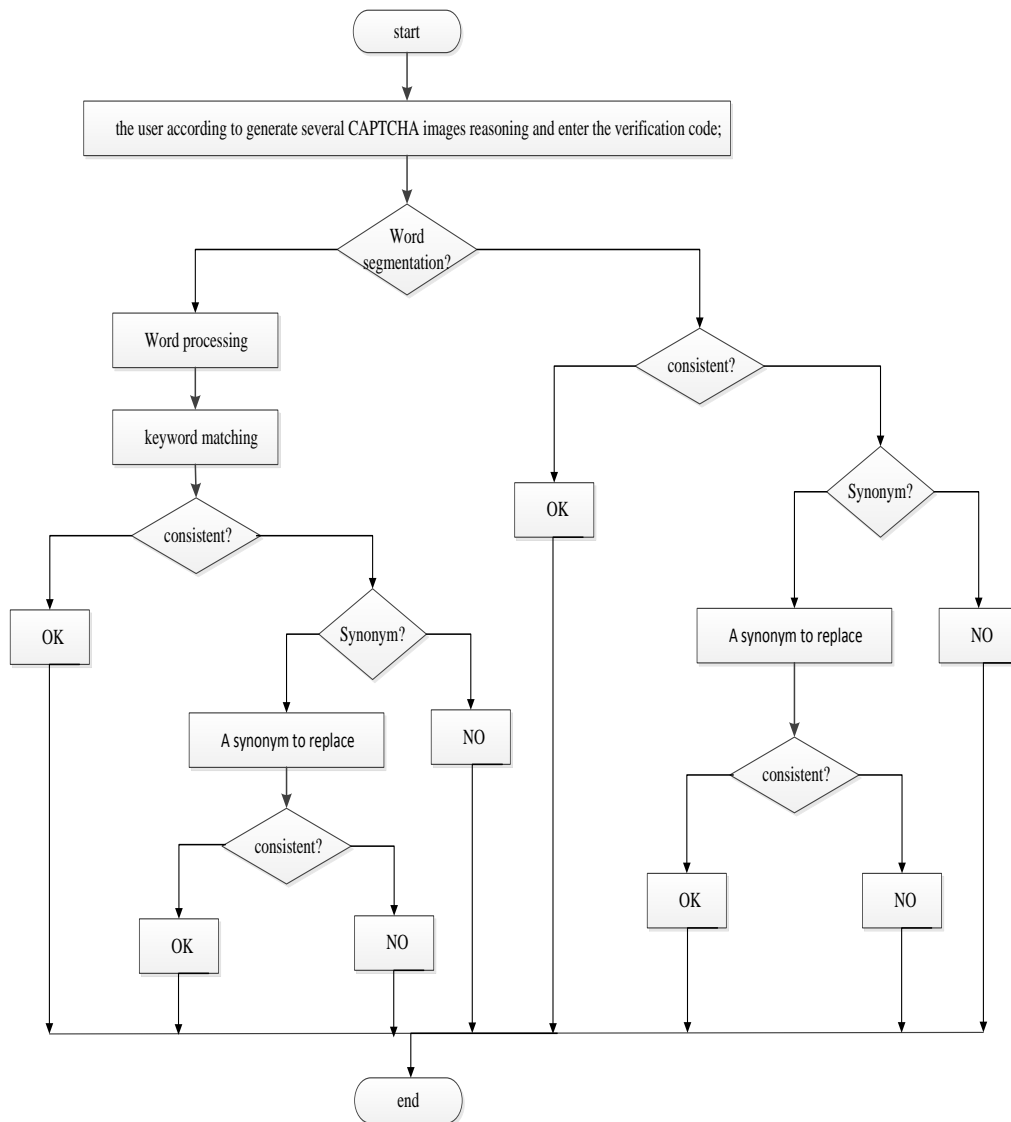


Figure 2. CAPTCHA validation algorithm

### 3.4 CAPTCHA implementation

The realization of the CAPTCHA is shown in Figure 3.

Username :

Password :

Can fly  
 Will lay eggs  
 The animal is :

[Forgot Password?](#)

Figure 3. Authentication code implementation

## Conclusion

The CAPTCHA has a variety of forms, but the development of CAPTCHA recognition technology causes a hidden danger for the security of the CAPTCHA. In order to improve the security of the CAPTCHA, a new kind of CAPTCHA based on knowledge base is put forward combining the implicit CAPTCHA, which is based on semantic information question and answer and the professional system. This new CAPTCHA can significantly improve the security of the professional system. The CAPTCHA designed in this paper is suitable for professional system but not for general system, such as E-mail.

## Acknowledgments

This work was supported by the scientific research fund project of China Agricultural University Yantai academy (YT201311, 201201Ja), the science and technology plan project of Beijing education committee (Grant no. KM201510853006) and key scientific research project of Beijing Polytechnic College (Grant no. bgzykyz201502, bgzykyz201503).

## References

- 1 JI Zhi-gang. "Principles and Prevention of HTTP Attacks Based on Identifying Code Recognition". *Computer Engineering*, 2006,32(20), pp. 170-172.
- 2 Ying Xiao-min. "The Research on User Modelling for Internet Personalized Services". PhD thesis, National University of Defense Technology, 2003.
- 3 Luis von Ahn, Manuel Blum, John Langford. "Telling Humans and Computers Apart Automatically". *COMMUNICATIONS OF THE ACM*, 2004, 47(2), pp. 57-60
- 4 Tao.R, Song Y.E, Wang Z.J. "Ambiguity function based on the linear canonical transform". *IET Signal Processing*, 2012, 6(6) , pp. 568-576.
- 5 WANG Bin-jun, WANG Jing-ya, DU Kai-xuan, etc. "Research on attach and strategy of CAPTCHA technology". *Application Research of Computers*, 2013, 30(9), pp. 2776-2779.
- 6 Mori G, Malik J. "Recognizing objects in adversarial clutter: breaking a visual CAPTCHA". *IEEE Conference on Computer Vision & Pattern Recognition (CVPR)*, 2003, 1, pp. 124-141.
- 7 ZUO Bao-he, SHI Xiao-ai, XIE Fang-yong, etc. "A Neural Network Based Approach to Recognizing the Verification Code". *Computer Engineering & Science*, 2009, 31(12), pp. 20-22.
- 8 Huang Sai-ping, Xu Ming. "Recognition and Improvement of Identifying Code". *Journal of Nanjing Normal University (Engineering and Technology Edition)*, 2009, 9(2), pp. 84-88.
- 9 LIU Jian-Wei, YAN Lu-Feng. "Comparative Study of Knowledge Representation". *Computer Systems & Applications*, 2010, 20(3), pp. 242-246.
- 10 Zhang Xuan-ping, Gao Hui, Zhao Zhong-meng. "The Rule Representation for Knowledge in Database Style". *Computer Engineering and Applications*, 2002, 38(1), pp. 200-202.
- 11 Zhang Qi-yu. "Research and Design of Spam Email Filter System based on Bayesian algorithm spam". M.S. Thesis, Qufu Normal University, 2006.

\*Corresponding author: Qiyu Zhang, Yantai Academy, China Agricultural University, No.2006, Coastal middle Road, Gaixin District, Yantai City, Shandong Province, China. E-mail: rcraingo@163.com