

# IMPLEMENTATION OF INDIVIDUAL CUSTOMIZED NEWS NOTIFICATION MODEL WITH WEB CRAWLING

<sup>1</sup>KYOUNG-ROCK CHUNG, <sup>2</sup>KOO-RACK PARK, <sup>3</sup>SEONG-HYUN PARK, <sup>4</sup>YOUNG-SUK CHUNG

<sup>1</sup> School of Computer Engineering, Kongju National University, Korea

<sup>2</sup> Key Laboratory of Computer Science & Engineering, Kongju National University

<sup>3</sup> School of Computer Engineering, Kongju National University, Korea

<sup>4</sup> School of Computer Engineering, Kongju National University, Korea

E-mail: <sup>1</sup>mtgood@naver.com, <sup>2</sup>ecgrpark@kongju.ac.kr, <sup>3</sup> darp1234@kongju.ac.kr, <sup>4</sup>merope@kongju.ac.kr

## ABSTRACT

In modern society, there are many news media that cannot be compared with the past because the news is made and reproduced by individual social media besides the official press. For example, newspapers were issued once a day to receive new news every day in the past. Now, with the development of information and communication technology, it is possible to check newspaper news in real time using PC and mobile internet, so that it can receive information anytime and anywhere quickly. In modern times, it is possible to see the news quickly through many media, but there are many articles that disappear without being aware of anything other than important events for a long time. If the event you are interested in is always in the news, you can keep track of all the events, but if a larger, more exciting new article is treated as a news story, you will not see an article about what happened afterwards. It is also missed. Many of the most recent events have been reported in so many news stories, but there are many cases in which the truth or progress of the event is not revealed as important as it was when the first event occurred.

In this paper, to solve these problems, we implemented a notification function that informs the news information of the individual keywords to the mobile application by constantly observing the news using the crawling technology.

Keywords: *NEWS, Database, Web crawling, Web scraping, Information retrieval.*

## 1. INTRODUCTION

With the development of internet, it has been possible to obtain a variety of information from web. The amount of data on the Web is growing at 40% per year, which is expected to increase from 4.4 ZB (zettabytes) in 2013 to 44.4 ZB (zettabytes) in 2020 [1]. There is a lot of information on the web, of which daily news is also large. In the past, most people received news on TV or newspapers, but now they are gaining such information from the wired Internet web or smartphone mobile apps. The era of seeing news in newspapers has been a long time ago, and it has only just begun to look at Internet news on personal computers, but now more people are getting the latest information via mobile. The way of contacting news is adapting quickly to the age according to the development of mobile device and communication technology, but the quality of news is not higher than before. In Korea, unlike other countries, most of the news is searched on the main

page of the portal site, and the article is focused on the exposed ranking information. As news agencies provide free articles to portal sites, the revenue of media companies is determined by the number of clicks or ads on search sites. So, it was overflowing with unnecessary articles because of duplicate registration of fishing articles and articles of irritating and sensational titles, which are more in - depth than high - quality articles and in - depth reports. As it became more important to edit the portal site in terms of content and level than the content and level of the article, the portal site operator has more power than the media company. Therefore, it takes a lot of effort and time to get the exact information that you want because of the editorial direction of each portal site and the news invasion of fake press.

This paper proposes customized news notification model based on web crawling. Web crawling is a process of extracting data from web sites. It uses hypertext transmission protocol to directly access World Wide Web and bring

information [2][3][4]. Web crawling is generally implemented with the use of web crawler. The URL and title of the news found using the keyword notifies by the mobile app.

This paper is comprised of as follows: chapter 2 describes relevant work; chapter 3 presents the proposal of customized news notification model based on web crawling; chapter 4 shows the implementation of the proposed model; chapter 5 describes conclusions and future search.

## 2. RELATED WORK

### 2.1. RSS

Typical ways to view news on the Internet or mobile is to use RSS provided by search sites or to use news categories provided by each portal site. RSS (Rich Site Summary; originally RDF Site Summary; often called Really Simple Syndication) is a type of web feed which allows users to access updates to online content in a standardized, computer-readable format. These feeds can, for example, allow a user to keep track of many different websites in a single news aggregator. The news aggregator will automatically check the RSS feed for new content, allowing the content to be automatically passed from website to website or from website to user. This passing of content is called web syndication. Websites usually use RSS feeds to publish frequently updated information, such as blog entries, news headlines, audio, video. An RSS document (called "feed", "web feed", or "channel") includes full or summarized text, and metadata, like publishing date and author's name. RSS (Rich Site Summary) is a contents expression method used in news or blog websites. A site administrator shows the contents of a website in the type of RSS [5].

There is the study on RSS as a marketing means for internet shopping malls [6]. A study was conducted to receive information without using a separate collection process using RSS to create a data broadcasting service process [7].

### 2.2 Web Crawler

A Web crawler (also known as a spider) is a Internet bots or software agents that systematically navigate the World Wide Web for the purpose of web indexing (web crawling or web spidering)[8][9]. Web search engines and some other sites use Web crawling or spidering software to update their web content or indices of others sites' web content. Web crawlers copy pages for processing by a search engine which indexes the

downloaded pages so users can search more efficiently. A Web crawler begins from the URL list called seeds. It recognizes all hyperlinks of a page to update a URL list, and recursively revisits the updated URL list [10][11][12].

There is a study that analyzes what kind of social networks are constituted among users by crawling web contents by answers of Naver KIN [13]. In order to efficiently search popular web pages, there is a research that applied the web crawling method to the graph search method [14].

### 2.3. HTML DOM

The DOM (Document Object Model) is a structure definition of all the elements constituting the web screen. It allows access to the internal structure of the elements on the screen. Using DOM allows asynchronous data access to dynamic views, which is very important for AJAX[15]. The DOM also allows you to dynamically change a web page and access and manipulate that page. The DOM allows you to integrate styles, values, and more in HTML. The model to manipulate the values of the screen content and its object is called DOM TREE. DOM NODE is one element of the tree structure, the most basic unit of DOM TREE [16] [17]. In this way, programs and scripts can efficiently access, update, replace, and delete documents.

Figure 1 is a DOM TREE structure diagram.

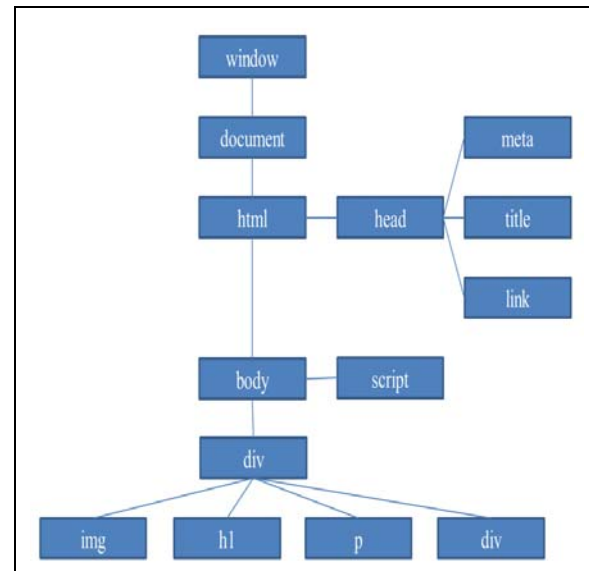


Figure 1: DOM TREE structure diagram

### 2.4 Multi-thread

A thread is a detailed execution unit that processes work in a process. A thread is a program

that allows a program to perform multiple tasks at the same time. Multi-threading allows a process to handle more than one task. That is, several threads operate in the same program to process the work at the same time, and it has the advantage that the processing speed of the software can be increased [18] [19].

**2.5 RESTful api**

RESTful api is an API that follows the REpresentational State Transfer architecture. (REST: means specifying a resource via an HTTP URL and defining its behavior through the HTTP Method (Post, Get, Put, Delete).) The actual definition is a web service architecture that conforms to ROA (Resource Oriented Architecture). It is intuitively "accessing objectized services using URI (Uniform Resource Identifier) and HTTP methods" and uses intuitive URIs. It is like accessing resources and manipulating them using HTTP methods. In addition, RESTful is a Web service middleware that allows service providers to respond directly to their resources without intermediate intermediaries when requesting resources such as graphics, images, text, audio, and hyperlinks based on ROA [20].

The following Figure 2 is a RESTful API structure.

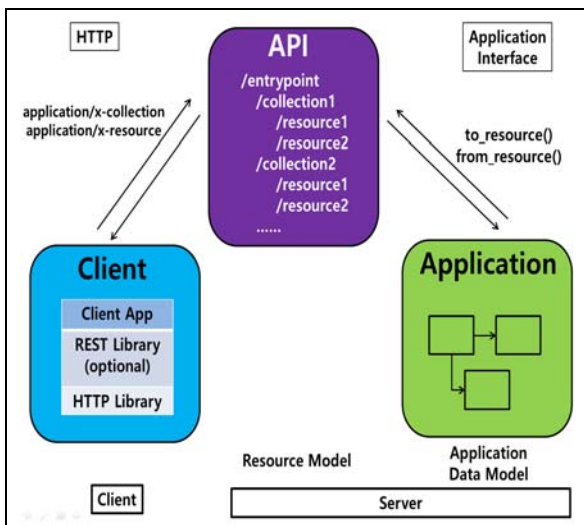


Figure 2: RESTful API structure

The main advantage of REST is that it is easy to use, and although there are some non-standard standards in common use, it is REST if only the features are kept well and the service is easily accessible. It is possible to implement REST without a framework or language-supporting

framework. This means that existing resources can be used as they are, so much effort is not required to introduce REST. REST was introduced in 2000 by Roy Fielding, one of the founders of the Web (HTTP). Because we believe that the current architecture does not use much of the original design excellence of the Web, we introduced a network-based architecture that can take full advantage of the Web. The elements of REST are largely composed of three elements such as resource, method, and message. One of the great features of REST is that it uses the existing web standard called HTTP, so it can use the existing infrastructure used on the web as it is. You can use the HTTP protocol-based load balancer or SSL, as well as caching, one of the most powerful features of HTTP. Considering that 60% to 80% of transactions in a typical service system are lookup transactions like select, caching HTTP resources on a web cache server can have many advantages in terms of capacity and performance. You can implement caching using the Last-Modified tag or E-Tag used in the HTTP protocol standard.

It is defined that the same purpose can be pursued based on the restriction condition, and the conditions are as follows.

1. Client / server architecture: Separate into a consistent interface.
2. Stateless: The client's context between each request should not be stored on the server.
3. Cache Handling Functions: As in the WWW, clients must be able to cache responses.
4. Tiering: Since there is no way to determine if a client is going through a target server or an intermediate server, it is possible to improve the scalability of the system by providing load balancing and shared cache of intermediate servers.
5. Code on demand: By providing a Java applet or JavaScript, the server can be extended by sending the logic that the client can execute.
6. Interface Consistency: Simplifies the architecture and breaks it down into smaller pieces, helping each part of the client and server to be distributed independently.

Based on these constraints, we separate tasks into front-end and back-end developers and perform tasks in accordance with the above conditions. You

can also speed up your work and avoid code twists with independent configuration. It can send and receive data in the desired type (XML, JSON, RSS) and is easy to link with various multiplatforms [21].

### 3. CUSTOMIZED NEWS NOTIFICATION MODEL BASED ON WEB CRAWLING

#### 3.1 Targeted News Process

Figure 3 is a flowchart for monitoring the user's news search keyword and target news for processing in this paper.

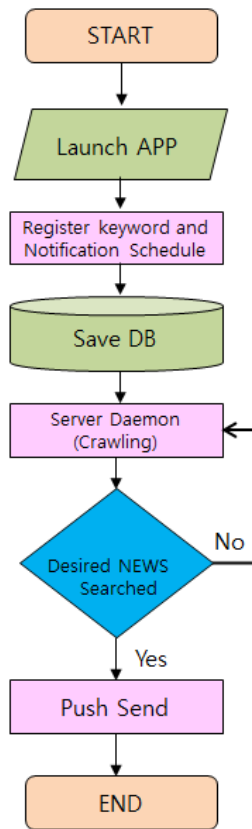


Figure 3: Targeted News Processing Flowchart

First step: Launch an app for user.

Second step: Register keywords to search for News and the schedule to notify the news you want.

Third step: Store the input data by app and use for monitoring crawling.

Fourth step: Run the crawl periodically on the server every 30 minutes.

Fifth step: Repeat the crawl every 30 minutes until the news extracted by the crawl matches the news desired by the user.

Sixth step: Send a push message to the user when the news wanted found.

#### 3.2 System Configuration

The following figure 4 is a system configuration diagram implemented in this paper.

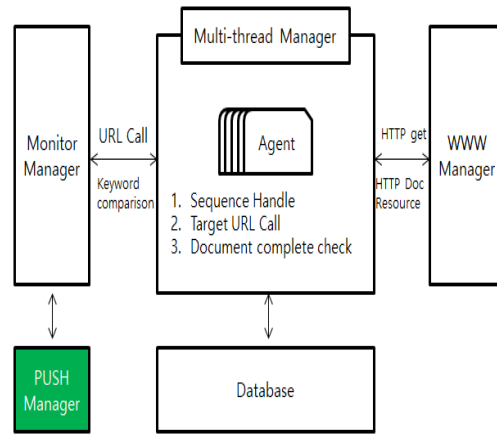


Figure 4: System configuration

The structure of the system implemented in this paper consists of four managers in total.

1. The Monitor Manager compares the keyword requested by the user and sends the message to the user's app by forwarding the event to the push manager when it is judged to be the desired news.

2. The Multi-Thread Manager configures and manages the threads so that they can extract data simultaneously to process multiple users' requests quickly. Then Agent accesses URL received from Monitor Manger, fetches the data of the site, and sends it to Monitor Manager again.

3. The WWW Manager manages and supervises the basic information of the web and delivers the HTML DOM and resources of the connected site to the Agent.

4. Push Manager delivers the message to the user registered with target URL and news title delivered from Monitor Manager by push using app.

### 3.3 Daemon Crawling

The following figure 5 is some of the daemon crawling sources running on the server.

```
keywordlist = db_keyword()
CallUrl("http://www.koreatimes.co.kr/") // call site
dat = CheckTagN("a")
web_data(dat)
Function web_data(data)
.
for(int i=0;i<=data.size();i++)
{
    if(keywordlist.indexOf(data[i])>-1)
    {
        data_txt=read_txt(data[i])
        AlramPush("skdiefkle.....");
    }
}
.
End Function
Function AlramPush(device_id)
.
if(device_id<>"" )
{
    Sender sender = new Sender("ekghdk.....");
    String regId = device_id;
    Message message = new
Message.Builder().addData("msg", alram_msg)
                .build();
    List<String> list = new
ArrayList<String>();
    list.add(regId);
    MulticastResult multiResult2;
    multiResult2=
sender.send(message, list, 5);
    if (multiResult2 != null) {
        List<Result> resultList
= multiResult2.getResults();
        for (Result result :
resultList) {
            System.out.println(result.getMessageId());
        }
    }
}
.
End Function
Function CallUrl(baseUrl)
.
driver_sebu = driver_sebu_m.getPage(baseUrl)
wait(driver_sebu, "")
```

```
.
End Function
Function CheckTag(tag,chkhtml)
.
List<HtmlElement> dic_atag = ((HtmlPage)
driver_sebu).getElementsByTagName(tag)
.
return dic_atag
End Function
```

Figure 5: Server daemon crawling

The implemented algorithm is designed to crawl news on the server every 30 minutes.

First: In the db\_keyword () function, the user fetches the search target keywords registered through the application as an array.

Second: The CallUrl () function accesses the site where users want to search for news.

Third: Take the HTML DOM of the searched site and get the title and news content of the news via the CheckTagN () function.

Fourth: In the HTML DOM obtained through the web\_data() function, Find the text by using the function of CheckTagN() and compare with the registered search target keyword.

Fifth: Compare keywords and send a push message to user's app if the same keyword exist. At this time, transmit only url and news title.

### 3.4 App Baseline Data Input Source For Users

Figure 6 below is a part of function that registers target keyword in user application.

```
public static void setup_insert()
{
    rs = st.executeQuery("select idx_no from
init_data where device_id="+device_id+"");
    if (rs.next())
    {
        String Query = "UPDATE init_data set
keyword1="+ keyword1 +", keyword2="+
keyword2 +", keyword3="+ keyword3 +",
        .... where
device_id="+device_id+"";
        Statement stmt = con.createStatement();
        rowCount = stmt.excuteUpdate(strQuery);
        //update count
    }
}
```

```

else
{
String Query = "INSERT INTO init_data
(device_id,keyword1,keyword2.....;
Statement stmt = conn.createStatement();
stmt.excuteUpdate(strQuery);
}
request_msg( "Save ok");
}
    
```

Figure 6: App baseline data input source for users

If the data registered in the application is json type, the above source code checks whether the existing data is registered through the executeQuery () function. If the registered data already exists in the corresponding user, the data is updated. if the data does not exist, Register as new one. When the update and insert are complete, the app sends a save completion notification.

**3.5 Json Communication**

The following figure 7 is a part of the function for data json communication with the app.

```

public static String user_senddata(String
keyword1, String keyword2,String keyword3,
String user_device.....) {
db_con();
JSONObject jsonObject = new
JSONObject();
JSONArray personArray = new
JSONArray();
JSONObject personInfo = new
JSONObject();
try {
String qu="SELECT * FROM
alam_master where device_id="+user_device+" ";
ResultSet rs =
statement.executeQuery(qu);
while ( rs.next() ) {
personInfo = new
JSONObject();
personInfo.put("item_cd",
rs.getString("item_cd"));
personInfo.put("url",
rs.getString("url"));
personInfo.put("title",
rs.getString("title"));
personInfo.put("text",rs.getString("te
xt"));
personInfo.put("open_yn",rs.getStrin
g("open_yn"));
personInfo.put("create_dt",rs.getStrin
g("create_dt"));
    
```

```

personArray.add(personInfo);
jsonObject.put("DATA",
personArray);
}
} catch (Exception e) {
}
return jsonObject.toString();
}
    
```

Figure 7: Json communication

The implemented algorithm is used as part of the source for communicating with the app and the server when receiving data from the app for the user.

First: The db\_con() function connects the DB.

Second: JSONObject creates data in string form as a json object.

Third: Returns the object converted to json.

**3.5 Json Communication**

The following figure 8 is the db layout for data collection.

DATABASE	MASTER	Table Name	INIT_DATA				
Table Space Name		Entity name					
Function							
NO.	Column name	Attribute name	Key	Null	Type	Length	Etc
1	idx_no			NO	int		
2	keyword1			YES	varchar	50	
3	keyword2			YES	varchar	50	
4	keyword3			YES	varchar	50	
5	due_date			YES	varchar	4	
6	device_id			YES	varchar	250	
7	create_dt			YES	datetime		
8							
9							
10							
11							
12							

Figure 8: Table Layout

This is a db layout for storing the data input by the user's application. The description of the column is as follows.

1. " idx\_no" is a unique number sequentially generated each time data is generated.

2. "keyword1" is the content and title of the news to be searched by the user, and is the first keyword among various ones.
3. "keyword2" is the content and title of the news to be searched by the user and the second keyword.
4. "keyword3" is the content and title of the news to be searched by the user and the second keyword.

When searching, it was configured to search from the title and contents of news included from keyword 1 to keyword 3.

5. Due\_date is configured to notify even if the same news is exposed due to redundancy, but only to inform about the news which is exposed after a predetermined time limit.
6. The device\_id is configured to store the unique key value of the device needed to send a push message to the mobile.
7. "create\_dt" is the date the user saved the content through the app.

#### 4. EXPERIMENTAL RESULT

The following figure 9 shows the screen for registering the news desired by the user in the user application.

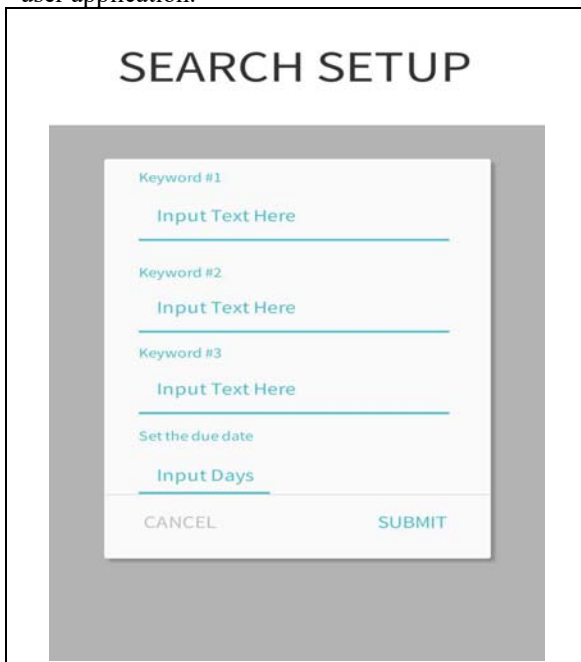


Figure 9: User Registration Screen

In order to test and run the implemented system, we set up to register three keywords that the user needs. Because it is often the case that the news of the matching keyword is duplicated, the notification date is registered in order to notify the user of the news after a certain period of time for the same keyword. In order to send the push to the user's app, it is not displayed on the screen, but the device id is acquired internally and the device id is also stored when the above data is stored.

The following figure 10 is a screen to notify the user when a push is received from the user application.

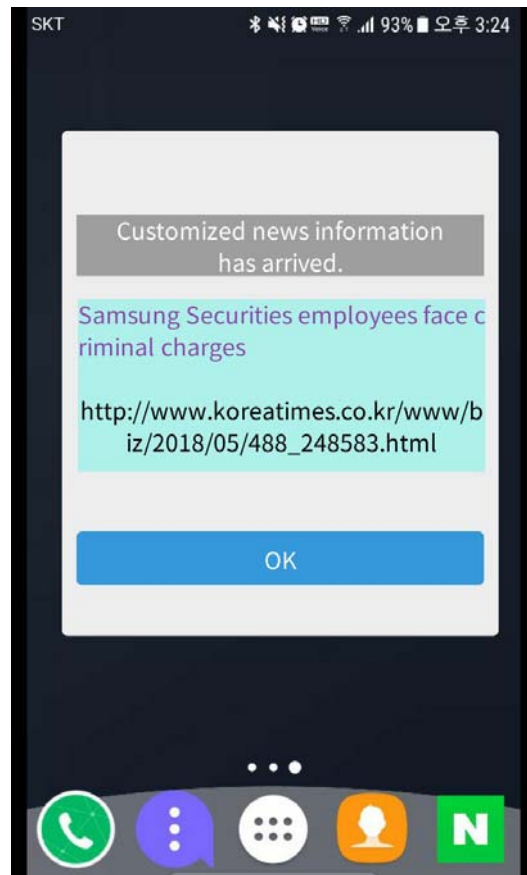


Figure 10: App push notification

The crawl daemon sends a push message to user's app when news that matches the user's registered keyword is detected by crawling. Then, the news title and url are exposed on the user's phone as shown above. When the user clicks the url in the above screen, the user can go to the news and read the article.

The following figure 11 is the screen that runs on the crawl daemon

```

===== start
==data & title=====Thu May 10 14:50:33
KST 2018 ---K-pop band INFINITE's Sungkyu to
join Army on May 14
Kim Sung-kyu, leader of the popular boy band
INFINITE, will begin his military duty next week,
his management .....
=====
----## Time required(sec .0f) : 0.001sec
=====
==data & title=====Thu May 10 14:50:34
KST 2018 ---May offers various events,
performances
=====
----## Time required(sec .0f) : 0.0sec
=====
==data & title=====Thu May 10 14:50:34
KST 2018 ---TVXQ powers through concert with
smash hits
By Kwak Yeon-soo
Much like Leonardo DiCaprio in "The Great
Gatsby," K-pop duo TVXQ threw a swanky party-
like concert, leading a .....
=====
----## Time required(sec .0f) : 0.003sec
=====
==data & title=====Thu May 10 14:50:35
KST 2018 ---'Burning' features beleaguered youth
of the time
By Park Jin-hai
"For the young generation of today, engulfed in
helplessness and rage, the world would come as a
big conundrum,.....
=====
----## Time required(sec .0f) : 0.001sec
=====
==data & title=====Thu May 10 14:50:35
KST 2018 ---2 Koreas resume cooperation in
films
NK actors may come to South for film fest
By Park Jin-hai
In October, a local film festival may feature North
Korean actors and directors.
Following the historic inter Korean Summit last
month, where the amicable atmosphere between
President Moon Jae-in .....
=====
----## Time required(sec .0f) : 0.001sec
=====
.
    
```

Figure 11: Crawl Daemon Launch Screen

The crawling daemon was developed with JAVA and compiled with jdk 1.7. I have brought the contents of the title and body of the crawled HTML DOM and indicated the total fetch time.

The imported data is implemented to judge whether it matches the user's request keyword and send a PUSH message.

The following figure 12 shows the JSON (JavaScript Object Notation) structure used in this system. We applied RESTful API method which is easy to apply for legacy interworking.

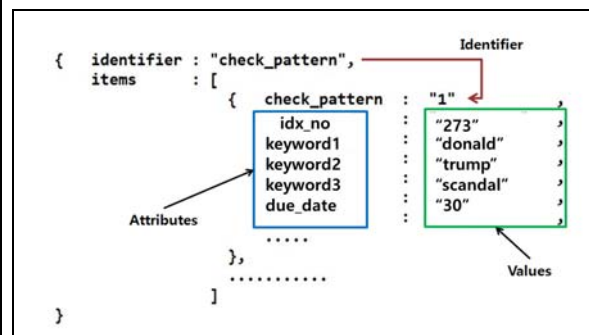


Figure 12: JSON Structure

JSON is a lightweight data interchange format that is easy for the human eye to recognize and is easy to analyze and generate in the system. JSON is a text format that can be directly used by a large number of programmers, not just a specific development language. These attributes make JSON an ideal data exchange language.

### 5. CONCLUSIONS

Until recently, newspapers and TV were used for news. Today, many people are watching the news on the internet and mobile. Various portal sites on the Internet offer news from various media companies, but the news providing directions for each site's news editorial did not escape from the one-sided method like old paper newspaper or TV. The freedom of information access is still low compared with the development of technology, so it is not easy to find news articles with desired information. And there are many cases where the news that user wants to go out after a long time. In this paper, we implemented web crawling system to find out the news that the user wants to know. As a result, users were able to see the news they wanted to know by app notification even after the user forgot.



Future research will continue to study to compare with each other the news to express a different point of view, even the same news.

#### REFERENCES:

- [1] Executive Summary Data Growth, Business Opportunities, and the IT Imperatives, <https://www.emc.com/leadership/digital-universe/2014iview/executive-summary.htm>
- [2] Junghoo Cho, "CRAWLING THE WEB: DISCOVERY AND MAINTENANCE OF LARGE-SCALE WEB DATA", stanford university(2001), pp.1-2
- [3] Kyoung-Rock Chung, Koo-Rock Park, Seong-Hyun Park, Young-Suk Chung. "A study on Notification Model of Specific News Information Using Web Scraping", International Workshop on Convergence Informaiton Technology, (2017) December21-23;Busan, korea
- [4] Dongmin Seo, Hanmin Jung, "Intelligent Web Crawler for Supporting Big Data Analysis Services", JOURNAL OF THE KOREA CONTENTS ASSOCIATION 13, 12(2013)
- [5] J.H.Jeon, S.Y.Lee, "Trend and Prospect of the Web 2.0 Technology, Electronics and Telecommunications Trends", ETRI(2006), Vol.21, No.5 pp.141-146.
- [6] Chae Hyuk Gi, Park Sangun, Kang Juyoung, "Applying RSS Marketing on Internet Shopping Malls Based on AISAS Model", The Journal of Society for e-Business Studies.13, 3 (2008)
- [7] Yunyoung Jang, Yoon-Chul Choy, Soon-Bum Lim, "Development of Data Broadcasting Service Creation System Based on Web 2.0 and RSS", KOREA INFORMATION SCIENCE SOCIETY.35, 1B(2008)
- [8] Wan-Sup Cho, Jeong-Eun Lee, Chi-Hwan Choi, "Refresh Cycle Optimization for Web Crawlers", JOURNAL OF THE KOREA CONTENTS ASSOCIATION.13, 6 (2013)
- [9] Hie-Cheol Kim, "Design and Implementation of a High Performance Web Crawler", Addison-Wesley, KOREA(2003)
- [10] Eun-Jeong Shin, Yi-Reun Kim, Jun-Seok Heo, Kyu-Young Whang, "Implementation of a Parallel Web Crawler for the Odysseus Large-Scale Search Engine", Journal of KIISE : Computing Practices and Letters.14, 6(2008)
- [11] Hye-Suk Kim, Na Han, Suk-Ja Lim, "Web Crawler Service Implementation for Information Retrieval based on Big Data Analysis", Journal of Digital Contents Society 18, 5(2017)
- [12] S.J.Kim, "A Comparative Study on Models of Web-based Information Seeking Behavior", Journal of the Korean Society for Information Management, 21, 2 (2004)
- [13] Chaeun Lee, Jinwook Jang, "Development of Social Data Collection System using Web Crawling", KCC2016, (2016) June 29-july;JeJu, Korea
- [14] Jinil Kim, YooJin Kwon, Jin Wook Kim, Sung-Ryul Kim, Kunsoo Park, "Effective Web Crawling Orderings from Graph Search Techniques", KOREA INFORMATION SCIENCE SOCIETY.
- [15] So-jin Nam, Do-Hoon Kim, Wan-Jung Kim, Yong-Hyuk Kim, "Dom-based Content Extraction for Improving Performance of Web Service", Korea Information Science Society,pp.218-222,Oct 2003.
- [16] Kim Kyuheon, Park JungWook, Kim Byungchul, "Effective Method to Change Multimedia Scene Configuration Information Using DOM Update", Journal of Broadcast Engineering, Volume 18(1), pp. 43-58, 2013.
- [17] Hyeon-Sook Lee, Jae-Hong Jeon, Doo-Kwon Baik, "A Study on the Generation of Dynamic Overview Map(DOM) in Hypertext", Korea Information Science Society, pp.1039-1042,Oct 1992.
- [18] Jung Hyun Im, "A Code Generation Tool for Multi-threaded Software Development", Addison-Wesley (Publisher Name), KOREA, 2015.
- [19] Jung Hyun Im, "Tuning the Thread Population in a Multi-threaded Software Development Process based on the DEVS formalism", Addison-Wesley (Publisher Name), KOREA, 2015.
- [20] Song, Byung-Kwen, 2014, "Performance Analysis of Web Service Middleware based on SOAP/RESTFUL," Journal of IKEEE, Vol. 18, No. 1, pp. 146~151.
- [21] Yong-Ju Lee, 2013, "Resource Matchmaking for RESTful Web Services," The Journal of Korean Institute of Information Technology, Vol. 11, No. 8, pp. 135~143.