# 全国地质资料馆数字地质资料馆主站日志数据集

1,2　　　1,2　　　1,2　　　1,2　　　1,2
1,2　　　1,2　　　1,2　　　1,2

1.　　　　　　　　　　　100037　2.　　　　　　　　100037

摘要：

IP

2014—2017

关键词：
数据服务系统网址　http://dcc.cgs.gov.cn

## 1　引言

14　　　　　　　　　　　2016

2016

2016　　　　　2012

6

2010

1987

E-mail　gxuezheng@mail.cgs.gov.cn

2016

2013

1

表 1　数据库（集）元数据简表

2014—2017

*.accdb

1.70 GB

http://dcc.cgs.gov.cn

"                                                          "

1212010004000150018

IP

.accdb                                                          .accdb

## 2　数据采集和处理方法

### 2.1　数据来源

IP

❶                                                          1

### 2.2　数据处理和应用

服务技术流程图



图 1　数字地质资料馆服务模式及数据获取流程

2

3

4



图 2　国内访问用户区域分布图

529.5

5 294 535 448.7 4 486 616

186.4

1 864 312

15

（万次）



图 3 国外访问用户区域分布图

24.6 246 333

4.7 47 091 15 4.6

45 927 18

表 2　国外访问国家关键词排序

| | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 1 | 1:20 | 1:25 | 1:50 | 1:20 | H45C003004 |
| 2 | 1:25 | 1:20 | 1:20 | | |
| 3 | | | | | |
| 4 | 1:20 | | | | 1:20 |
| 5 | 1:20 | oilfields | | | |
| 6 | 1:25 | | | | |
| 7 | 1:25 | H4820 | 1:20 | 1:20 | |
| 8 | 1:20 | 1:20 | | 1:50 | 1:25 |
| 9 | | 1:20 | | 1:50 | |
| 10 | | | | | 1:25 |

2　　　　　　　　　IP　　　　　"　1　20　　　""　1　25　　　　"

1　20

1　50　　　　　　　　　IP



图 4　用户检索关键词统计图

2017

## 3 数据样本描述

Access

Access                                                    IP                                3

                                                    4

IP                              2014—2017                              IP

                        2014—2017                                    IP

表 3　数字地质资料馆网站访问 IP 地址记录数据表

| | | |
|---|---|---|
| 1 | | 2293B1F0579C7019E05341015A0A617B |
| 2 | IP | 14.215.222.217 |
| 3 | | /Data/FileList.aspx?MetaId=E928A0F55D2F7A73E0430100007F3D67&type=zw |
| 4 | | 631795983@qq.com |
| 5 | | |
| 6 | / | 2015−10−21 9:21 |

表 4　数字地质资料馆网站搜索关键词记录数据表

| | | |
|---|---|---|
| 1 | IP | 59.71.224.2 |
| 2 | | |
| 3 | / | 2015−10−17 19:16 |
| 4 | | 224C06EB72A00920E05341015A0A506A |
| 5 | | |
| 6 | | |

## 4 数据质量控制和评估

                                                    603              6 034 025

                2017      3      17      18:40      3      18      20:10   2017      10      5      14:37      10

7      12:06

                        IP

## 5 访客分析

　　　　　　　　　　　　　　1　20　　　　　　52 859　　　1　20　　　　　　　　33 290

　1　25　　　　　21 583　　1　50　　　　　　11 813　　　　　　3 392　　"

"　　　　　　　　　　　　　　　　　　　　"　　　"　　　　　　97.2%

　　　　　　1　20　　　　　　　　　　　　　　　　　　43.0%

　　　7.2%

## 6 结论

　1

　2

　3

　4

致谢：

Esri

注释：

## 参考文献

，　　　，　　　，　　　，　　　，　　　　　. 2016.
　　[J].　　　　，25(2): 73−76.

　　　，　　　，　　　，　　　，　　　，　　　　. 2016.　　　　　　　　　　　　　　　　　[J].
　　，25(2): 92−96.

　　　，　　　，　　　，　　　　. 2016.　　　　　　　　　　——
[J].　　　　，43(2): 691−697.

　　. 2016.　　　　　　　　　　　　　　　　[J].　　　　，18: 31−33.

　　，　　　. 2013.　　　　　　　　　　　[J].　　　　，8: 79−85.

　　. 2010.　　　　　　　　　——　　　　　　　　　[J].　　　　，
5: 62−68.

# The Dataset of User's Accessing Log to Digital Geological Library of National Geological Archives of China

GAO Xuezheng[1,2] LI Chenyang[1,2] WU Xuan[1,2] KONG Zhaoyu[1,2]
QI Fanyu[1,2] SHANG Yuntao[1,2] JIA Liqiong[1,2] LI Xiaolei[1,2] GUO Hui[1,2]

(1.*Development and Research Center of China Geological Survey*, Beijing 100037, *China*; 2.*National Geological Archives of China*, Beijing 100037, *China*)

**Abstract:** In order to accurately find out the needs of users of geologic data and eliminate the information gap between data users and data managers, the National Geological Archives of China (NGAC) has conducted the data collection of user's accessing logs to the master station of the Digital Geological Library of NGAC. The automatic recording method is applied to register the user's location, their searching keywords, IP address and other related information for constructing the accessing log dataset. In order to make better use of these accessing data, a standardized processing method and quality control system are adopted. The dataset provides the accessing records to the Digital Geological Library of NGAC's master station from the year 2014 to 2017, which realistically reflects user's behaving habits while obtaining geological data. It also provides a sound basis for the further constructions of geological data service website, the development and utilization of geological data, and the geological data management and services.
**Key words:** Digital Geological Library; website; user's accessing log
**Data service system URL:** http://dcc.cgs.gov.cn

## 1 Introduction

The NGAC is the largest and the most complete professional archives in China's national geological industry. Currently, more than 140, 000 kinds of geological data have been collected (Wang Xinchun et al., 2016), covering such fields as regional geological survey, mineral prospecting and exploration, marine geological exploration, geophysical, geochemical and remote sensing geological survey, hydrogeology, engineering geology and environmental geological survey, theoretical geological researches, geotechnical researches (Gao Xuezheng et al., 2016). Stepping upon the on-site lending-borrowing working flow, the Digital Geological Library (DGL) of NGAC is a key national-level geological information infrastructure that features itself with digitalization. It reforms the traditional workflows with digital approaches, and provides the worldwide academic

society with the up-to-date on-line geological information services (Kong Zhaoyu et al., 2016). The construction had been launched since the year 2012. After 6 years of endeavors, the DGL can provide multiple services to the academic communities, including geological data index services, metadata services, geological map services, thematic services, and inter-links and reciprocal checking services.

During the operation of the DGL, a large amount of user's accessing log information is recorded, which can help better understand user's needs, accessing habits, and their awareness and acceptance about the website, and capture more information so as to guarantee safe and stable operation of website (Zhao Guohong, 2010). It also enables managers to accurately grasp the developing tendencies of the website, figure out more specific ways to develop the website in the future (Xie Xiaoping, 2016); meanwhile, it will also help improve the quality of geological data services, improve higher recognition and satisfaction worldwide, and further consequently better solve a series of problems such as dull and dry unitary data servicing way, low efficiency, high cost of acquisition, insufficient utilization of electronic geological data, and many other constraints towards the geological data development as well. Al these efforts will guide the website construction towards the directions of more intelligence, more precise, and more initiative (Yu Shi et al., 2013).

Table 1 lists the metadata of the accessing log to the master station of Digital Geological Library of NGAC.

**Table 1   Metadata table of Database (Dataset)**

| Items | Description |
|---|---|
| Database (dataset) name | The dataset of user's accessing log to Digital Geological Library of National Geological Archives of China |
| Database (dataset) authors | Gao Xuezheng, Development and Research Center, China Geological Survey (NGAC)<br>Li Chenyang, Development and Research Center, China Geological Survey (NGAC)<br>Wu Xuan, Development and Research Center, China Geological Survey (NGAC)<br>Kong Zhaoyu, Development and Research Center, China Geological Survey (NGAC)<br>Shang Yuntao, Development and Research Center, China Geological Survey (NGAC)<br>Qi Fanyu, Development and Research Center, China Geological Survey (NGAC)<br>Jia Liqiong, Development and Research Center, China Geological Survey (NGAC)<br>Li Xiaolei, Development and Research Center, China Geological Survey (NGAC)<br>Guo Hui, Development and Research Center, China Geological Survey (NGAC) |
| Data acquision time | 2014—2017 |
| Geographic area | Nationwide |
| Data format | .accdb |
| Data size | 1.70 GB |
| Data service system URL | http://dcc.cgs.gov.cn |

Continued table 1

| Items | Description |
|---|---|
| Fund project | Geological Exploration Project of China Geological Survey "National Geological Data Integration and Data Collation" (121201004000150018) |
| Language | Chinese |
| Database (dataset) Composition | The dataset is composed of two parts, namely: Digital Geological Library's website accessing IP address registration table .accdb; Digital Geological Library's website keywords searching record registration table.accdb |

## 2   Data Acquisition and Processing Methods

### 2.1   Data Source

Accessing log data of master station of Digital Geological Library of NGAC is from the backend database in the Digital Geological Library website, recorded in real-time throughout the day. The log contents include all the visiting users' login sites, search queries, browsing information, thematic queries, geological map accessing and a series of other related processes. The dataset records in detail the IP addresses of the accessing user, the keywords searched, the time of accessing, the type of user, the location of the user and other related information; and it can also distinguish whether the user refers to the map information or the document information.

Fig. 1 shows Digital Geological Library of NGAC service model and data acquisition process.



Fig. 1    Shows Digital Geological Library of NGAC service model and data acquisition process

## 2.2 Data Processing and Application

Data recording, data extraction, warehousing and other system processing is conducted for accessing log dataset of master station of Digital Geological Library of NGAC. Firstly, the website access user's information and operation are recorded; secondly, the data entries are extracted for the data items through the structured data statement; and thirdly, they are stored into the database according to a unified data format.

Since the data in text format can not directly show the changes of the content in each field of the dataset, the data samples are displayed through the domestic (Fig.2) and foreign (Fig.3) accessing user's distribution area, and user's keyword searching statistics (Fig.4).



Fig. 2    Domestic accessing user's geographical locations

Based on the statistics on the geographical distribution of domestic accessing users in this dataset, the total times of user's accessing reaches 5.295 million times (5, 294, 535 times); if that from the internal LAN is omitted, it would be 4.487 million (4, 486, 616 times). Beijing, Shanghai, Guangzhou are the top three cities in terms of accessing times, with a total of 1, 864, 000 times (1, 864, 312 times). Fig. 2 shows that the users in the Yangtze River Delta and Pearl River Delta are more concentrated and that from the western China are much less.

The statistical data show that there is a relatively higher times of access to and search in the website of Digital Geological Library of NGAC by users from domestic cities of China such as Shanghai, Dongguan and Hangzhou, where traditionally less geological activities have been conducted. On the contrary, cities in traditional geological major provinces rank low in access times, for example, Qingdao in Shandong Province only ranks the 15th in China. This is mainly due to the following reasons. First, propaganda for network service provided by the Digital Geological Library of NGAC is insufficient, so many people engaged in geological activities have no idea about the Digital Geological Library of NGAC and the data published by it, resulting little use of the geological data

（Million times）



Fig. 3    Foreign accessing user's geographical locations

by traditional major geological provinces. Second, the first tier cities of China, Yangtze River Delta region and Peal River Delta region are leading in information technology and more sensitive to network, resulting in a higher accessing times from these regions. Third, the Digital Geological Library of NGAC mainly provides search and page view of geological data, no download function is available at present, and users can not acquire the data they are interested in. This is unfavorable for utilization of the digital geological data. Furthermore, the need of geological activities, particularly field geological activities to utilize digital geological data relies on Internet accessibility. However, at present field geological activities in China do not have the condition to access the Internet in most cases.

Among the foreign countries, the accessing times to NGAC from the United States reaches 246, 000 (246, 333 times), ranking the 4th. That from Germany is 47, 000 (47, 091 times), ranking the 15th. Ukraine has 46, 000 visits (45, 927 times) to NGAC, ranking the 18th. That from other countries such as Britain, France, Russia, Netherlands, Japan, Singapore, Canada, Italy, UAE, Sweden, Brazil, Finland, Australia and other countries are also respectively more than one thousand times. Therefore, we can see that foreign countries are also very concerned about the content released by the website of Digital Geological Library of NGAC.

**Table 2    Keyword sorting of foreign country**

| Ranking | Country | Key words 1 | Key words 2 | Key words 3 | Key words 4 | Key words 5 |
|---|---|---|---|---|---|---|
| 1 | United States | 1:200, 000 geological maps | 1:250, 000 geological maps | 1:500, 000 geological maps | 1:200, 000 hydrogeological maps | H45C003004 |
| 2 | Germany | 1:250, 000 geological maps | 1:200, 000 hydrogeological maps | 1:200, 000 geological maps | Laishui | Distribution map of groundwater resources in China |

Continued table 2

| Ranking | Country | Key words 1 | Key words 2 | Key words 3 | Key words 4 | Key words 5 |
|---------|---------|-------------|-------------|-------------|-------------|-------------|
| 3 | Ukraine | nothing | nothing | nothing | nothing | nothing |
| 4 | France | 1:200, 000 geological maps | Xianju | Jiangyong | Qiqihar | 1:200, 000 gravity data distribution map in China |
| 5 | United Kingdom | 1:200, 000 geological maps | oilfields | geological maps | Daging | Hydrology |
| 6 | Russia | 1:250, 000 geological maps | Ordos Basin | Ordos | Hebei | Shanxi |
| 7 | Netherlands | 1:250, 000 geological maps | H4820 | 1:200, 000 geological maps | 1:200, 000 hydrogeological maps | China |
| 8 | Japan | 1:200, 000 hydrogeological maps | 1:200, 000 geological maps | tungsten ore | 1:500, 000 geological maps | 1:250, 000 geological maps |
| 9 | Singapore | Bauxite | 1:200, 000 geological maps | Guangxi | 1:500, 000 geological maps | Aluminum mine |
| 10 | Canada | Xi'an | Shanxi | groundwater | Jiapigou | 1:250, 000 geological maps |

Table 2 shows that '1:200, 000 geological maps' and '1:250, 000 geological maps' have been the two most commonest keywords for search by access from foreign IP addresses. Other keywords concentrated by foreign access users include '1:200, 000 hydrogeological maps', '1:500, 000 geology map' and 'Shanx'. Abnormalities exist in access by IP addresses in Ukraine, where no search keyword is generated but there is concentrated continuous accessing.

According to the statistics of keywords being searched in the dataset, the top 5 keywords are '1:200, 000 geological maps' (52, 859 times), '1:200, 000 hydrogeological maps' (33, 290 times), '1:250, 000 geological maps' (21, 583 times), '1:500, 000 geology map' (11, 813 times), Xinjiang (3, 392 times). 'Geological map' has become the main content of user's searching targets. Among the top 5, the retrieval ratio of 'geological map' reaches 97.2%, being the absolutely main target keyword; while that of '1:200, 000 geological maps' occupies 43% in the top 5, and is 7.2% in all the keywords, being the most frequently retrieved keywords by users.

It is worth to note that Hebei, which is not a geological major or giant province at all, has ranked the 3rd in being retrieved. According to analysis, we consider that this is critically related with the decision of the Central Committee of Communist Party of China and the State Council of China to establish Xiongan New Area, Hebei Province. After establishment of this new area, ministries of the State Council began to prepare New Area Planning, China Geological Survey carried out field geological survey in Xiongan New Area. All these activities increased users' concern about Hebei Province and keyword search regarding Hebei Province. This shows that users' retrieving behavior and need are closely related with government policies and current hot affairs. In our future work, it is advisable to provide high quality service products to meet the specific needs of users.
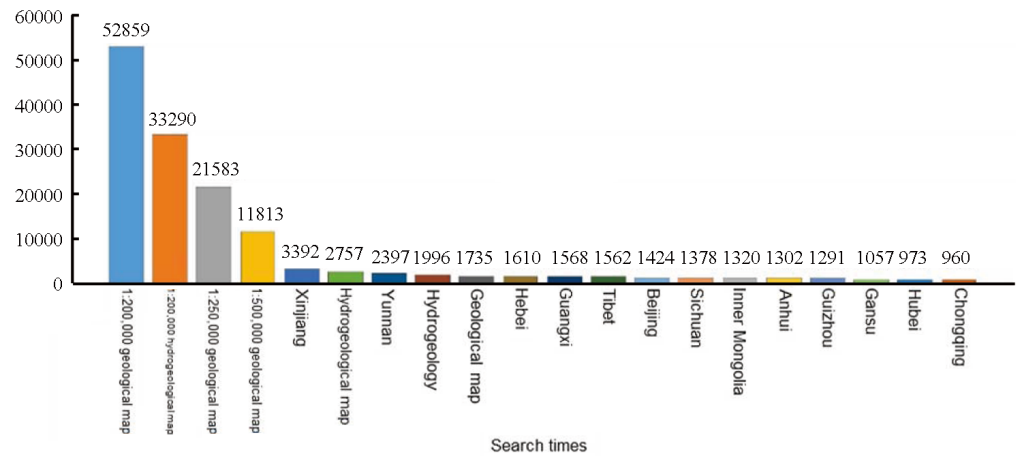
Fig. 4　User's retrieving keywords statistics

## 3　Data sample description

The Dataset of User's Accessing Log to DGL of NGAC is of Access format database, including two Access format database files, (1) Digital Geological Library's website accessing IP address registration database, and (2) Digital Geological Library's website keywords searching record database. Among them, the first one records the Digital Geological Library's visitors IP address, accessing data content, whether or not the real name registered users, user's accessing geographical location and accessing time during the years 2014—2017; while the second one records the Digital Geological Library's visitors IP address, searched keywords, search time and other data items during the years 2014—2017.

Table 3　Digital Geological Library's website accessing IP address registration database

| No. | Field Name | Format | Example |
|-----|-----------|--------|---------|
| 1 | Only string | Short text | 2293B1F0579C7019E05341015A0A617B |
| 2 | IP address | Short text | 14.215.222.217 |
| 3 | Access path | Long text | /Data/FileList.aspx?MetaId=E928A0F55D2F7A73E0430100007F3D67&type=zw |
| 4 | User account | Long text | 631795983@qq.com |
| 5 | User location | Long text | China, Guangdong, Foshan |
| 6 | Access time | Data/Time | 2015−10−21 9:21 |

Table 4　Digital Geological Library's website keywords searching record database

| No. | Field Name | Format | Example |
|-----|-----------|--------|---------|
| 1 | IP address | Short text | 59.71.224.2 |
| 2 | Keywords | Long text | Zhejiang Huzhou |
| 3 | Access time | Data/Time | 2015−10−17 19:16 |
| 4 | Only string | Short text | 224C06EB72A00920E05341015A0A506A |
| 5 | User account | Short text | Anonymous User |
| 6 | User location | Long text | China, Hubei, China University of Geosciences (Wuhan) |

## 4   Data Quality Control and Evaluation

This user's accessing log dataset of master station of Digital Geological Library of NGAC contains all the data since the digital data logger started to record the data, with two data tables in the dataset recording a total of 603 million pieces (6, 034, 025) of data. Due to the twice power-supply upgrading maintenance, the master station's accessing log records are not continuous, respectively, from 18:40 on March 17, 2017 to 20:17 on March 18, 2017, and from 14:37 on October 5 to 12:06 on October 7. However, the incomplete records only take a very small percentage of the total records, and do not affect the completeness, reliability, applicability and accuracy of the entire dataset.

The specific working process of data collection is as follows: starting from the very beginning of the user's accessing to the website, the data logger begins to record the user's IP address, accessing location, accessing content and other information, and meanwhile preprocess the recorded data according to the instruction to check whether there is abnormal recording or non-expected performs. According to the time sequence, accessing log data is stored to ensure the accuracy and effectiveness of warehousing. In order to ensure the safe and stable operation of the database, a regular database updating and maintenance strategy is formulated and strictly implemented, and the contents of the database are periodically evaluated. In the process of data extraction and analysis, the data items are unified based on format to ensure the overall consistency of the data content. In order to further improve the data quality, the project team will make more efforts in improving the mechanisms of dataset updating and maintenance, increasing the volumes of new information, and enhancing the versioning of accessing log data.

## 5   Analysis on Visitors

Statistical analysis on domestic and foreign users' search keywords has revealed that the top 5 search keywords are '1:200, 000 geological maps' (52, 859 times), '1:200, 000 hydrogeological maps' (33, 290 times), '1:250, 000 geological maps' (21, 583 times), '1:500, 000 geology map' (11, 813 times) and 'Xinjiang' (3, 392 times). 'Geological map' has become the main content searched by users. The search ratio of 'geological map' reaches 97.2%, being a dominant keyword in the top 5. The keyword '1:200, 000 geological maps' accounts for 43% search times of the top 5 keywords and 7.2% search times of all keywords, being the most frequently searched keyword.

In terms of geological visitors, the first tier cities of China and developed cities in Yangtze River Delta region and Peal River Delta region have a higher visitor volume, while western regions of China have a lower visit volume. In general, geological visitors of China are normal in access behavior, while many foreign visitors have abnormal access behavior. Analysis on search keywords, registration information, search language and search frequency has revealed that these search behaviors are mostly considered as web redirection to foreign server to hide real IP address by domestic agents. Therefore, among foreign visitors accessing the website, suspects of malicious request, theft and capture of information are not excluded.

## 6   Conclusions

(1) The Digital Geological Library has used the information system to comprehensively transform and upgrade the traditional geological archives. The paper-back predominated

processing pattern has been updated to the digitalization predominated pattern, the transformation of which has realized the technical innovations in China's geological data management and service. Consequently, the operation mode and organization style of the national and provincial geological data archives have both experienced vital changes.

(2) The user's accessing log dataset of master station of Digital Geological Library of NGAC is of a database that has digitally recorded the user's access to the archived digital resources. The relevant records can be used to track the whole process of user's access to the geological data network. Through statistical and correlation analysis, the multiple internal relations among such many factors as accessing users, accessing time and accessing content can be detected; the user's needs and searching habits can hence been better digitally supported.

(3) According to the analysis of the dataset content and the user's interested fields, the thematic data service is offered, and meanwhile the website home pages is newly arranged; all these efforts have further upgraded the service level of the Digital Geological Library of NGAC, reaching the purposes of precise, efficient and quick services.

(4) The content of the dataset, the operation patterns, and data analysis samples can provide references for the socialized service of the other digital archive agencies; it is especially valuable in promoting the qualified services of China's geological data.

**Notes:**

❶ Development Research Center, China Geological Survey. 2013. NGAC Core digital system R & D results report [R].

## References

Gao Xuezheng, Kong Zhaoyu, Qi Fanyu, et al. 2016. Research into the development of geological data collection and service in the National Geological Archives of China [J]. China Mining Magazine, 25(S2): 73−76 (in Chinese with English abstract).

Kong Zhaoyu, Shang Yuntao, Gao Xuezheng, et al. 2016. Research on construction of national geological data center [J]. China Mining Magazine. 25(S2): 92−96 (in Chinese with English abstract).

Wang Xinchun, Qi Fanyu, Li Xiaolei, et al. 2016. Research on the geological data integration and service: A case study of geological work in the equipped exploration area [J]. Geology in China, 43(2): 691−697 (in Chinese with English abstract).

Xie Xiaoping. 2016. Analysis of the Application of Web Visiting Statistics in the Archival User Research [J]. Lantai World, 18:31−33 (in Chinese with English abstract).

Yu Shiyang, Wang Jiandong. 2013. Government website analysis enters big data era [J]. E−Government, 8:79−85(in Chinese).

Zhao Guohong. 2010. Analysis of traffic impact factors of government portal websites: based on the comparison between Chinese and Japanese government websites [J]. E−Government, 5:62−68 (in Chinese).