# Research on Construction of Aluminum Industry Market Information Database

Jing Zhao[1], Yi Huang[1], Daling Li[2]
[1]Institute of Scientific and Technical Information of China
No. 15 Fuxing Road, Haidian District
Beijing, 100038, China
zhaojing@istic.ac.cn

[2]Tianjin Institute of Scientific and Technical Information
No. 22 Wujiayao Street, Hexi District
Tianjin, 300074, China

ABSTRACT. *In the process of economic globalization, the competitiveness of enterprises depends largely on the overall strength and innovation ability of industrial cluster. The development and construction of subject database becomes a major issue we confront, especially under the environment of the information age. This research takes market and product information of aluminum industry as subject, collects macro dynamics, resource development, and enterprises relevant to aluminum industry, and develops the aluminum industry market information database (ALIMID). With many subject database merged into an integrated enterprise resource database, we can provide resources and support services for innovation activities and promote the development of innovation and scientific research.*
**Keywords:** aluminum; database construction; specialized database;

1. **Introduction.** In the process of economic globalization, the competitiveness of enterprises depends largely on the overall strength and innovation ability of industrial cluster [1-2]. The development and construction of subject database becomes a major issue we confront, especially under the environment of the information age. This research takes market and product information of aluminum industry as subject, collects macro dynamics, resource development, and enterprises relevant to aluminum industry, including the status of the aluminum industry of world's major developed countries and regions, market prospects, resources reserves, exploitation and production information, as well as the

development of medium-sized enterprises, and develops the aluminum industry market information database (ALIMID). With many subject database merged into an integrated enterprise resource database, we can provide resources and support services for innovation activities and promote the development of innovation and scientific research. [3-5]

The rest of this paper is organized as follows. Section 2 introduces idea and framework. Sections 3 describes ALIMID knowledge system. Sections 4 describes ALIMID-Metadata design specifications. Sections 5 describes the logical framework and data organization of ALIMID. Sections 6 describes Data collection and annotation of ALIMID. Finally, a brief conclusion is given in Section 7.

2. **Idea and Framework.** Centering on enterprises' demands for information resources and knowledge, with the guidance of overall framework and norms of knowledge organization, this research builds aluminum industry market information database. The overall framework can be divided into four steps: firstly, define the metadata by making reference to the relevant national standards, international standards and industry standards; secondly, design the logical framework for database, build an information source standard as well as methods for data analysis and retrieval, and collect the raw data with related technologies; thirdly, make full use of existing editing technologies to clean, review and annotate the raw data to achieve precision control of the data; and finally, establish update and maintenance mechanisms to ensure database availability.[6-7] The technical roadmap is shown as Figure 1.
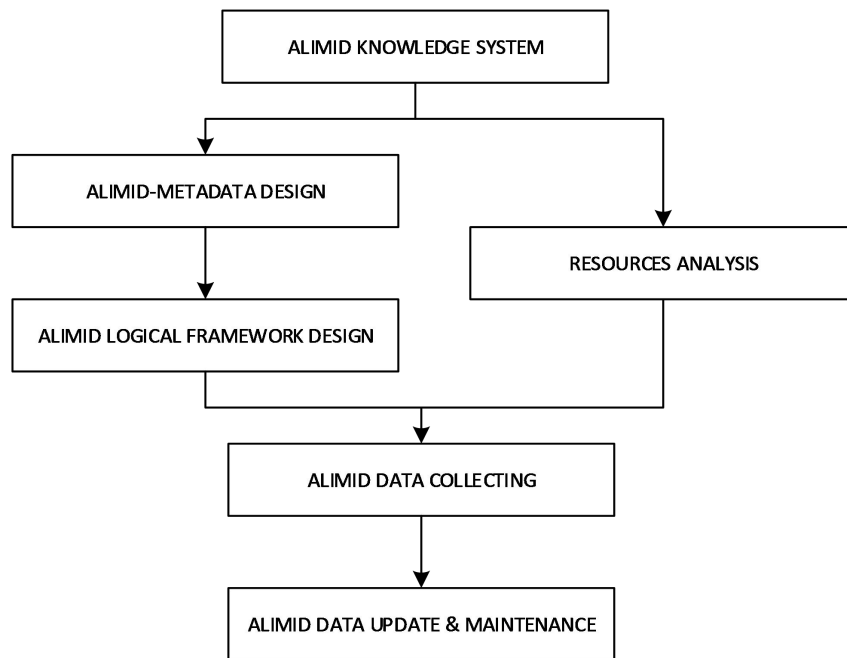


FIGURE 1. THE TECHNICAL ROADMAP OF ALIMID CONSTRUCTION

3. **ALIMID knowledge system.**

3.1. **Design principle and method of construction.**

**(1) Design Principles**

The construction of knowledge system means the integration of definition and affiliation relations of professional knowledge in the content level to archive the basic classification and navigation of resources effectively, as well as enhancing the integrity and comprehension of industrial knowledge.

According to the general requirements of ALIMID construction and the related national standards, international standards, industry standards, market information in the field of aluminum industry classification structure, we analyze and construct knowledge system of the aluminum industry market information, to achieve the top-level organization and indexing of aluminum industry market information data resources.

**(2) Construction Methods**

In the process of knowledge system construction, through a top-down view, we break the large complex problem confronted down into small simple ones, find the key to every problem by means of quantitative and qualitative analysis, and generate the aluminum industry market information concept tree.

Based on the aluminum industry market concept tree, we follow the compiling method of GB/T 13745 "subject classification and code" to build the system. In this knowledge system, the categories under the same hierarchical level form a parallel relation, and those between different levels form a superior-subordinate relation.

3.2. **Classification system.** According to the construction requirements for the database, based on the actual situation of domestic and foreign markets, we classify the data into two categories, i.e. industry information, supply and demand information. Combined with the discipline classification of categories and subcategories, 2 categories, 8 major categories and 19 subcategories are being constructed as shown in Table 1.

TABLE 1. CLASSIFICATION OF ALIMID

| CATEGORIES | PRIMARY CATEGORIES | MINOR CATEGORIES |
|---|---|---|
| INDUSTRY INFORMATION | NEW | DOMESTIC NEWS |
| | | INTERNATIONAL NEWS |
| | MARKET INFORMATION | MARKET ANALYSIS |
| | | ALUMINUM PRICES |
| | ENTERPRISE INFORMATION | MANUFACTURERS |
| | | MERCHANT |
| | | SERVICE PROVIDER |
| | CONFERENCE INFORMATION | CONFERENCE REPORTS |
| | | CONFERENCE |
| | POLICIES AND REGULATIONS | INDUSTRY POLICY |
| | | TARIFFS |
| SUPPLY AND | TENEMENTS POSTINGS | MINING RIGHT |

35

| DEMAND INFORMATION | | PROSPECTING TRANSFER |
|---|---|---|
| | MINING MACHINERY SUPPLY AND DEMAND INFORMATION | MINING MACHINERY SUPPLY |
| | | MINING MACHINERY DEMAND |
| | ALUMINUM PRODUCTS TRADE LEADS | ALUMINIUM INFORMATION |
| | | ALUMINUM INFORMATION |

By using linear classification, we encode the aluminum industry market information classification codes with six digits, which can be reduced or expanded depending on the circumstances. The information in this specification is mainly in the aluminum industry market, and we made a modest expansion as shown below:

TABLE 2. CATEGORIES AND CODES IN ALIMID

| CODE | CATEGORIES |
|---|---|
| 01 | INDUSTRY INFORMATION |
| 02 | SUPPLY AND DEMAND INFORMATION |
| 0101 | NEW |
| 0102 | MARKET INFORMATION |
| 0103 | ENTERPRISE INFORMATION |
| 0104 | CONFERENCE INFORMATION |
| 0105 | POLICIES AND REGULATIONS |
| 0201 | TENEMENTS POSTINGS |
| 0202 | MINING MACHINERY SUPPLY AND DEMAND INFORMATION |
| 0203 | ALUMINUM PRODUCTS TRADE LEADS |
| 010101 | DOMESTIC NEWS |
| 010102 | INTERNATIONAL NEWS |
| 010201 | MARKET ANALYSIS |
| 010202 | ALUMINUM PRICES |
| 010301 | MANUFACTURERS |
| 010302 | MERCHANT |
| 010303 | SERVICE PROVIDER |
| 010401 | CONFERENCE INFORMATION |
| 010402 | CONFERENCE |
| 010501 | INDUSTRY POLICY |
| 010502 | TARIFFS |
| 020101 | MINING RIGHT |
| 020102 | PROSPECTING TRANSFER |

| 020201 | MINING MACHINERY SUPPLY |
|--------|-------------------------|
| 020202 | MINING MACHINERY DEMAND |
| 020301 | ALUMINIUM INFORMATION |
| 020302 | ALUMINUM INFORMATION |
| 020303 | INDUSTRY INFORMATION |
| 020304 | SUPPLY AND DEMAND INFORMATION |

4. **ALIMID-Metadata design specifications.**

4.1. **Design process of ALIMID-Metadata.** Aluminum industry market information database metadata (ALIMID-Metadata) provides basic information about the aluminum industry market, which is the basis for the construction of topics in ALIMID.

It describes the aluminum industry market information based on content. It describes the structure of the aluminum industry market information, which can be used as the aluminum industry market information sharing and exchange infrastructure.

In ALIMID-Metadata, we define unique identifier for each and every entry in ALIMID, so as to provide the basis of uniquely identifies in aluminum industry.

According to ALIMID-Metadata structure, we can also make some expansion.

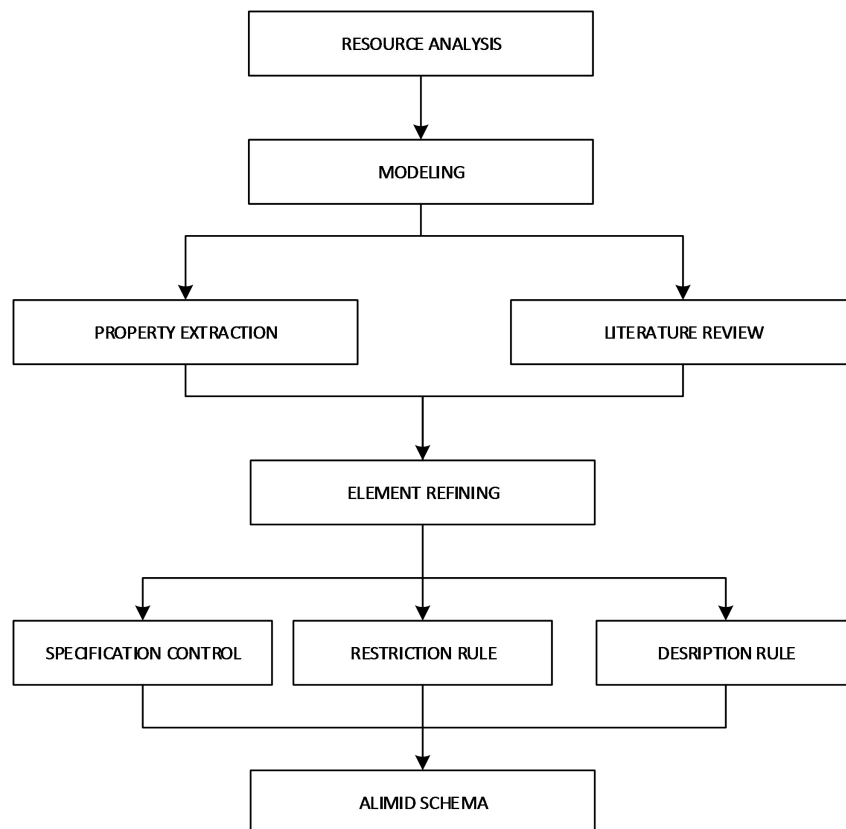The design process of ALIMID-Metadata is shown as follow:



FIGURE 2. THE DESIGN PROCESS OF ALIMID-METADATA

4.2. **ALIMI-Metadata set.** Based on common characteristics of the aluminum industry market information resources, ALIMID-Metadata set is a standard that can be used to describe the exploration, mining, production, processing and other aspects of market information resource.

ALIMID-Metadata set consists of two parts: descriptive information, metadata reference and contact information.

**(1) Descriptive information**

Descriptive information records the basic information of the database, such as the data set name, profile, creator, source, etc.

Each element has a descriptive label for users to read, and a unique ID for computer to process.

We use English words in lowercase in naming all the elements to make it easy for computer to mark and encode, as well as to ensure consistency with other languages; and we use Chinese in naming all the tags to facilitate users to read. Tag is just a semantic attribute of an element, designed to highlight the aluminum industry market information resources and metadata and to better reflect the name of the element in a particular semantics information within the original definition in specific applications.

TABLE 3. ALIMID-METADATA SET

| ELEMENT NAME | DEFINITION | NOTE |
| --- | --- | --- |
| ID | The unique identifier assigned to the aluminum industry market information | It refers to the aluminum industry market information unique identifier by taking a string compliant labeling systems identify and use the global unified coding rules. |
| Title | The name assigned to the piece of information | It makes resource well representation with the formal name. |
| Creator | The main responsibility of the entity providing the information content | it includes individual, organization and institution authors. |
| Keyword | Topic content of the resource | Keywords to describe the main content of the resource, with keywords or classification codes. |
| Publisher | Information on which entity bears responsibility to issue | Posted by including publisher individuals, organizations or institutions. It should be representative of the name used to identify the publisher of the entity. |
| PublishDate | publish date and time of the information | It uses ISO 8601 format, namely YYYY-MM-DD form. Wherein, YYYY is the calendar year, MM is the month of the year, DD is the day of the month. Example: 2003-04-01 represents April 1, 2003. |
| Source | original source to obtain | the original source that can be obtained as a |

| | the information | whole or as a part. |
|---|---|---|
| Language | Language of the information | It use RFC1766 format consisting of two characters (from ISO639) components. Such as "en" for English, "fr" for French. |
| ISSN | reference of the information, ISSN | A variety of content types and serials (such as newspapers, periodicals, yearbooks, etc.) have a unique identification codes of the assigned carrier type. |
| Text | full text of the information | It is the actual content of the original information in ALIMID. |

### (2) Metadata reference and contact information

Metadata reference and contact information records the information of metadata itself, including metadata standard name, metadata creation date, metadata contact information, etc. With this information, on one hand, a user can connect data collection information related to metadata records; on the other hand, metadata reference metadata information will facilitate maintenance personnel to modify and maintain the metadata.

TABLE 4. ALIMID-METADATA REFERENCE INFORMATION

| ELEMENT NAME | DEFINITION | NOTE |
|---|---|---|
| ContactName | Contact name associated with the database | When the contact is individual, this element should be filled in the name of the contact. |
| OrganizationName | Contact Unit | When the contact is a unit or organization, this element should be filled in the full name of this unit or organization. |
| Address | Contact address and zip code | Contact address and zip code to zip code, specific to the street, house number, box number or the unit, department name of an organization. |
| Fax | Contact Fax Number | Use "(area code) number" or "(Code) switchboard number - extension number" format. If there is more than one fax number, use a semicolon between numbers as a separator. |
| Phone | Contact telephone number | Use "(area code) number" or "(Code) switchboard number - extension number" format. If there is more than one phone number, use a semicolon between numbers as a separator. |

| Email | Contact e-mail address | Email address format, such as user@abc.com. If there is more than one email, use a semicolon between emails as a separator. |
|---|---|---|

## 5. Logical framework and database organization of ALIMID.

5.1. **Database Structure.** Aluminum industry market information database contains eight specialized database and an index database. Each specialized database makes an integration of knowledge system with common data. ALIMID then uses a unified view of the data mapping between the index database and specialized database, so as to realize the joint retrieval and integration of application in a deeper level.
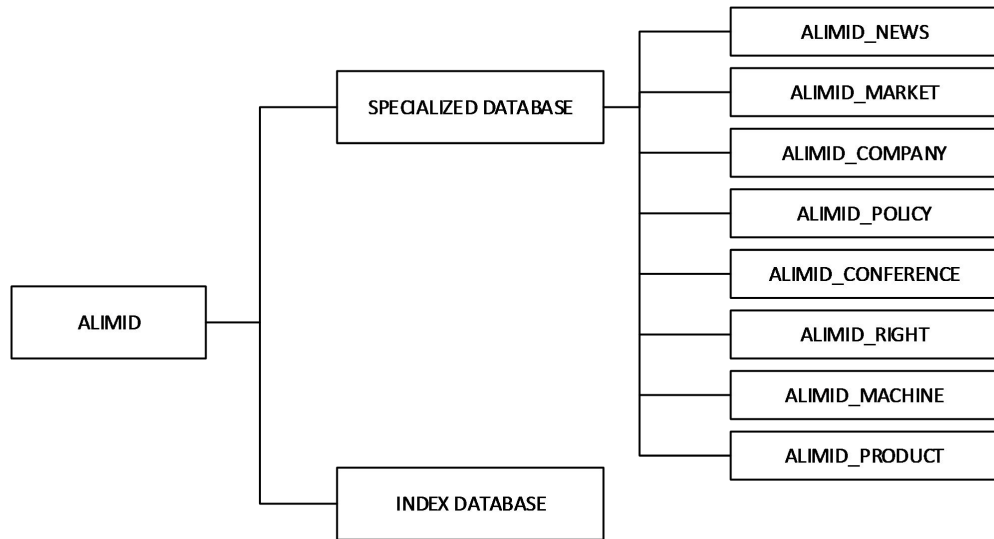


FIGURE 3. LOGICAL FRAMEWORK OF ALIMID

5.2. **Design of Index database.** The mapping information of common data model between specialized database, logical database and MADAI Metadata are used to generate an index database automatically. Index database contains three fields, namely, identification, classification, and relational database table.

TABLE 5. INDEX DATABASE TABLE DESIGN

| NO. | FIELD NAME | TYPE | LENGTH | DEFAULT | ALLOW NULL VALUE | PK |
|---|---|---|---|---|---|---|
| 1 | Id | int | 4 | | | √ |
| 2 | TypeCode | int | 4 | | √ | |
| 3 | Dbconname | varchar | 50 | | √ | |

5.3. **Design of Specialized database.** According to the eight categories in the knowledge structure, eight specialized database are constructed, namely ALIMID_News (aluminum news database), ALIMID_Markets (aluminum market database), ALIMID_Company (aluminum company database), ALIMID_Policy (aluminum policy database), ALIMID_Conference (aluminum conference database), ALIMID_Right (Aluminum right

40

database), ALIMID_machine (aluminum machine database), ALIMID_Product (aluminum product database).

5.4. **Data organization of ALIMID.** With diverse sources of information, the data reliability and credibility varies. The reliability and quality of information is the most important issue to build ALIMID. We establish a data resource survey and selection method, and determine a reasonable range of the collection of information sources and types of information included in the form, in order to ensure the integrity and authority of the collected information.

**(1) Data sources**

Data sources can be divided into two parts, data from the Internet and data from NSTL specialized database.

Data from Internet comes from four types of sites: the government website (People's Republic of China Ministry of Land website, Xinhua, China Economic Net, People's Daily, the state-owned Assets Management Committee, etc.); aluminum industry portal and related industry sites (China aluminum Website, China Aluminum Network, Global network of mineral rights, mineral resources network, mining machinery, China mining machinery Information Network, aluminum prices website, China Nonferrous website, the Yangtze River Nonferrous Metals, etc.); monetary and financial websites (CICC network, China Metal News, Reuters, Chinese commercial data center, the huge influx of information, Shanghai Nonferrous Metals website, China Yangtze River Nonferrous Metals Spot Metals website, etc.); local specialties website (China - ASEAN mineral resource net, Henan Nonferrous Metals, Nonferrous Metals Shandong, Anhui Province, land and mineral market network, etc.).

TABLE 6. SPECIALIZED DATABASE DESCRIPTION

| DATABASE NAME | SIZE | UPDATE FREQUENCY | TYPE |
|---|---|---|---|
| CNKI China Academic Journals Full-text Database | 29052131 items | Monthly | Journal |
| VIP scientific and technical journals database | 90G | Monthly | Journal |
| WANFANG DATA Dissertations of China | 2084383 items | Monthly | Thesis |
| CNKI China Masters' Theses Full-text Database | 1634162 items | Monthly | Thesis |
| CNKI China Doctoral Dissertations Full-text Database | 200614 items | Monthly | Thesis |
| MacroChina government decision support database | 100G | Monthly | Standard literature |
| SRIT Database | 150G | Daily | Standard literature |
| Dialog | Synchronized with | Synchronized with | Standard literature |

| | server database | server database | |
|---|---|---|---|
| Kompass International Enterprise Products Directory Database | Synchronized with server database | Synchronized with server database | Enterprise Products |
| Chinese enterprises and Product DATABASE | 124474 items | Bimonthly | Enterprise Products |
| National laws and regulations database | 145646 items | Half-yearly | Laws and regulations |

**(2) Selection Method**

A. Keyword section

We choose those words that can express the characteristics of the substantive contents of the database vocabulary, such as "aluminum", "bauxite" and so on; also those words that can reflect the search features and other features of the non-theme words, such as "author", "time", "type", etc. After the selection of keywords, we also check the full name and abbreviation, spelling form of those words, to find out the hidden theme concept with good truncation symbols, to avoid error and detection.

B. Query construction

We use logical operators to construct the action range of each term, utilize operators to adjust the position relationship between search terms, and use parentheses to establish search priorities, so as to achieve narrowing search range and improve the precision rate.

C. Query validation and refinement

After using query to conduct a search, we verify those results and choose those resources we need.

**6. Data collection and annotation of ALIMID.** We follow certain technical standards and specifications to process, handle, and import heterogeneous information resources, and integrate data from different sources to achieve data retrieval, statistical analysis, data mining and knowledge discovery.

The data collection methods we used include Internet data collection, documents collection, and relational database information collection.

6.1. **Internet data collection.** We use iRMS to download, monitor and store information that meets the requirements from a variety of dynamic and static sites and professional literature databases.
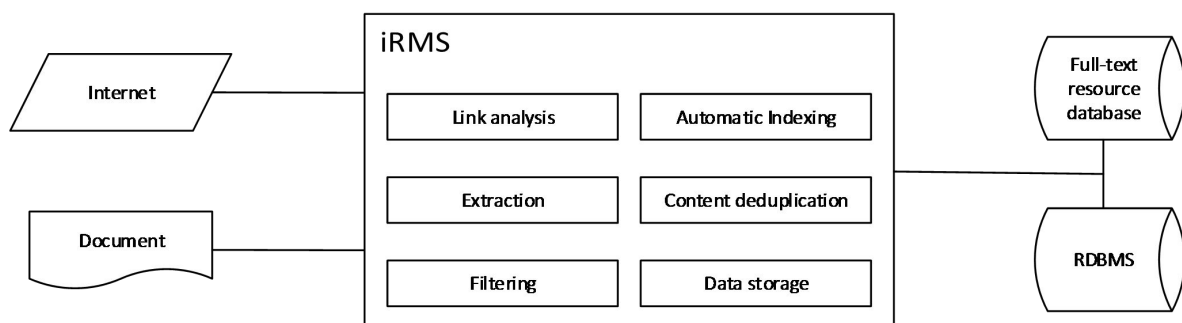


FIGURE 4. FRAMEWORK OF INTERNET INFORMATION COLLECTION

Its main features include:

Webpage downloading: automatically download Internet information, and support web-based templates for page metadata extraction and content extraction.

Database resources downloading: support metadata and full-text attachments downloading from network database resource, and details template-based resource for resource metadata extraction.

Automatic information filtering and classification: implementing the Internet for gathering information or document title based on content filtering, automatic classification set support, support link row of heavy weight and content.

Automatic Content Extraction and Indexing: based on customized template to achieve metadata extraction and content extraction, and automatic indexing.

Information export: support customized data output, and those date can be loaded into the full-text resource database or a relational database.

Real-time monitoring and automatic monitoring: through task scheduling setting, it can regularly and automatically track changes, automatically identify additional webpages, and documents so as to do real-time monitoring of Internet and documents.
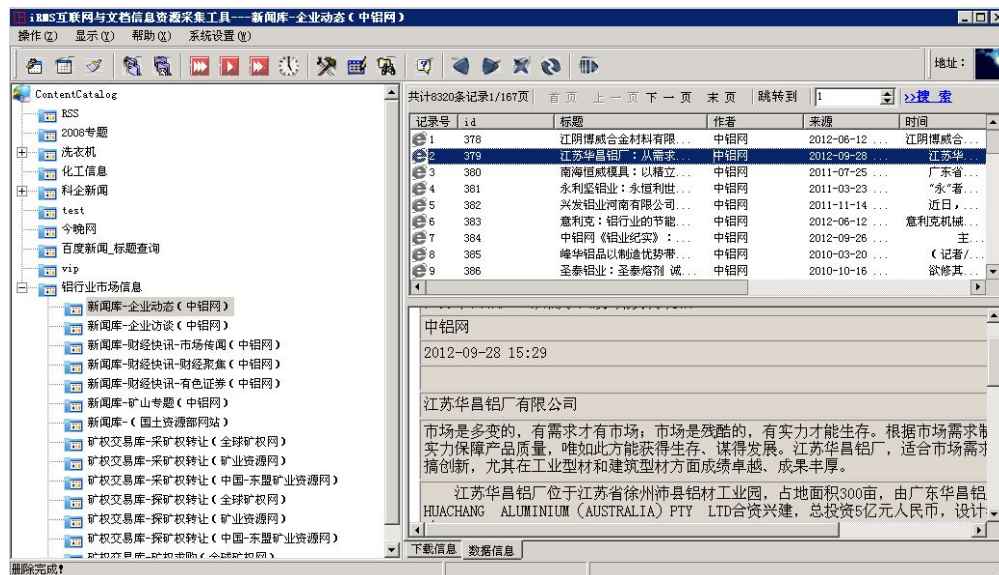


FIGURE 5. DATA COLLECTION CONFIGURATION AND RESULTS

6.2. **Document information collection.** We also use iRMS to collect electronic document information, mainly to complete documents in a variety of formats for catalog indexing, automatic scanning, text extraction, original link generation, automatic storage and other functions. For documents stored inside the organization server as well as personal computers, we can also use iRMS to collect, process, store and index documents.

It includes the following features:

(1) To organize documents by category: catalog classification hierarchy can be customized, scan documents in the root directory and automatically extract the directory as a document classification and indexing items.

(2) Metadata extraction: for some standard format documents, we can extract some characteristic data values, such as title, author, organization, abstract and other items information as metadata indexing with original document links automatic generation.

(3) Feature extraction: support TXT, HTM / HTML, RTF, WORD, EXCEL, PPT, PDF and other document feature extraction and text extraction; for graphics, images and other multimedia data, it support classification and indexing of multimedia information and implemented for loading, management and services.

(4) Parameter configuration: it can be run as a background document collection tools with configuration parameter, and it support mission planning and operation of automatic scheduling, data files directory monitoring and continuous updating of the document data.

**6.3. Relational database information collection.** For those data in business system that uses relational database as underlying storage, we can use a relational database synchronization tools DEC to archive data transformation and data storage confirming to specifications.

Its main features include:

(1) It can batch convert and integrate major relational database resources to the full text resource database system with task schedule in multiple ways so as to collect, extract, convert, integrate, synchronize data in relational database of government agencies and enterprises, to provide a unified resource integration and service application.

(2) With dual card, soft gatekeeper technology to ensure safe internal business application systems.

(3) In the data conversion process, we can specify the correspondence relation between fields, capable of processing large text, image streaming data, and exchange data in XML format.

6.4. **ALIMID annotation.**
6.4.1. **Data publishing.** In the data annotation phase, we use TRS Web content management system (TRSWCM) to process, edit and index the raw data we collected in previous steps.

In accordance with the table structure and metadata of specialized database, we import those xml files and other data sources into the system, then the existing data is displayed through the static publishing engine to provide basic data for editing and publishing in the next phrase.

With static pages generation technology, the contents of the database are being generated in HTML format to improve access efficiency. Before publishing, it provides preview function of home page, column, page, etc., and it XML formatted page output.

Data publish can be conducted in two ways: incremental publishing and complete publishing: (1) incremental publishing, to publish only the most recent qualifying manuscripts, and (2) complete publishing, to publish the entire column or even republish the entire contents of the website. We can specify only publish the data of the day, in which state the manuscript are allowed to publish, and the order of records in overview page.

FIGURE 6. TRS WEB CONTENT MANAGEMENT SYSTEM

**6.4.2. Data annotation.** With visual editing features of TRSWCM, we can edit, review, delete and annotate the data content via WYSIWYG editing, and the annotation process includes:

(1) Refining concept theme

Based on the title, abstract, introduction, conclusion, secondary title, charts and references information of each entry, we determine the theme of the entry, in case of any doubt we can browse the full text to further define the theme.

(2) Subject headings annotation

There is no need to make the standardization process to reveal the subject matter of any entry. Since subject headings annotation method is simple, and we only need to choose words that are approximate in natural language.

(3) Keyword extraction priority

We first extract keywords from the header of each entry, then from abstract when title cannot meet the requirements, but also browse the full text when necessary to conduct the extraction;

(4) Keywords number

The number of Keywords is generally 3-8, for keywords less than 3 cannot fully describe the theme of an entry, and keywords more 9 would broaden the range of an entry and cause information redundancy.

6.5. **ALIMID data collection results.** We conduct ALIMID construction according to the framework and process shown in previous, and we collect 11606 items on Aluminum Market Industry Database shown as Table 7.

45

TABLE 7. ALIMID DATA COLLECTION RESULTS

| Type | Total | Type | Total |
|---|---|---|---|
| Domestic News | 1267 | Mining Right | 28 |
| International News | 1057 | Industry Policy | 302 |
| Market Information | 431 | Tariffs Information | 265 |
| Domestic Enterprise | 441 | Conference Information | 1168 |
| International Enterprise | 844 | Aluminum information | 5743 |
| Machine Supply & Demand | 60 | | |

7. **Conclusion.** In the process of economic globalization, the competitiveness of enterprises depends largely on the overall strength and innovation ability of industrial cluster. The development and construction of subject database becomes a major issue we confront, especially under the environment of the information age. In view of this, this paper takes market and product information of aluminum industry as subject, collects macro dynamics, resource development, and enterprises relevant to aluminum industry, and develop the aluminum industry market information database (ALIMID). We conduct ALIMID construction according to the framework and process proposed, and we collect 11606 items on Aluminum Market Industry Database to provide resources and support services for innovation activities and promote the development of innovation and scientific research.

**REFERENCES**

[1]  Xu, Guo-Dong, Hong Ao, and Yuan-Guan She. "Current status and development trend of aluminum industry in world and strategy suggestions in China under background of sustainable development." The Chinese Journal of Nonferrous Metals 22.7 (2012): 2040-2051.

[2]  QIAO, XIAODONG, YAO LIU, and WEI JIN. "DESIGN AND PRACTICE OF SERVICE-ORIENTED PLATFORM FOR ENTERPRISE TECHNOLOGY INNOVATION." ICIC express letters. Part B, Applications: an international journal of research and surveys 6.4 (2015): 1219-1224.

[3]  LIU, YAO, YI HUANG, and YAN WANG. "RESEARCH ON THE KEY TECHNOLOGIES OF PYRIOS KNOWLEDGE SERVICE PLATFORM." ICIC express letters. Part B, Applications: an international journal of research and surveys 6.5 (2015): 1323-1328.

[4]  Li, Jie, Ruijia Wang, and Yao Liu. "Research on Semantic Metadata Online Auxiliary Construction Platform and Key Technologies." In Chinese Lexical Semantics, pp. 702-716. Springer Berlin Heidelberg, 2013.

[5]  Liu Yao, Haiqing Shi, and Deju Zheng. "Study on Semantic Annotation for Professional Literature." ICIC Express Letters. Part B, Applications: an International Journal of Research and Surveys 5, no. 5 (2014): 1383-1389.

[6]  QIAO, XIAODONG, MINGCHANG WANG, and ZHIJUN GUO. "RESEARCH ON ONTOLOGY CONSTRUCTION FOR PRODUCTION PROCESS OF ALUMINUM ELECTROLYSIS." ICIC express letters. Part B, Applications: an international journal of research and surveys 5.1 (2014): 257-264.

[7]  Guowu Nong, Xiaodong Qiao, Lijun Zhu, Yao Liu. On Construction and Implementation of Aluminum Enterprise Knowledge Management System. Digital Library Forum. 2013(010):37-43.