

DeepShovel: An Online Collaborative Platform for Data Extraction in Geoscience Literature with AI Assistance

SHAO ZHANG, Shanghai Jiao Tong University, China

YUTING JIA, Shanghai Jiao Tong University, China

HUI XU, Shanghai Jiao Tong University, China

YING WEN, Shanghai Jiao Tong University, China

DAKUO WANG, IBM Research, United States

XINBING WANG, Shanghai Jiao Tong University, China

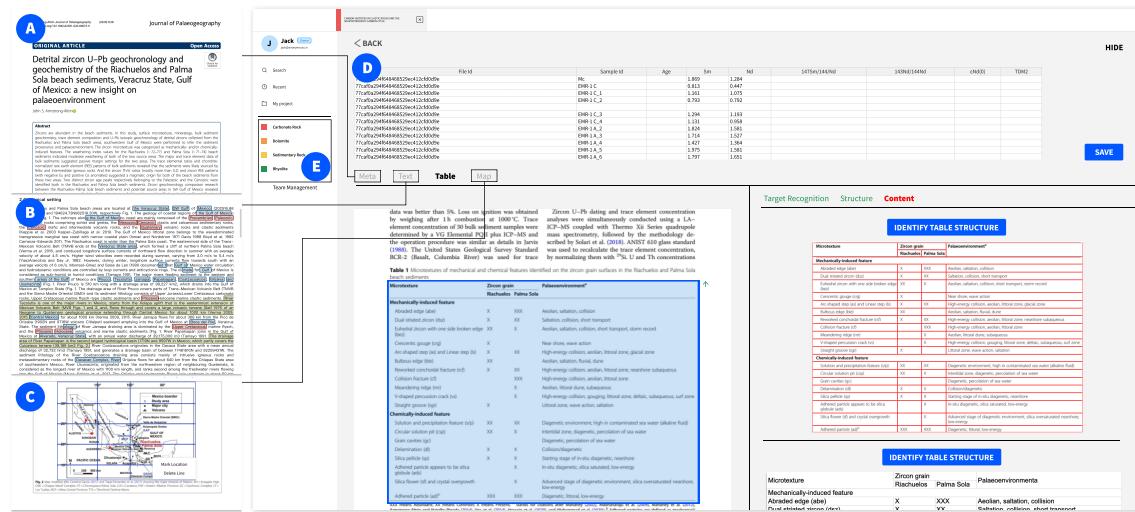


Fig. 1. User Interface of DeepShovel, an online collaborative platform for data extraction in geoscience literature with AI assistance. The main part of this figure illustrates the table extraction and integration functions (D), while the system can also support meta information extraction (A), text extraction (B), map recognition and location extraction (C), and team and document management (E).

Geoscientists, as well as researchers in many fields, need to read a huge amount of literature to locate, extract, and aggregate relevant results and data to enable future research or to build a scientific database, but there is no existing system to support this use case well. In this paper, based on the findings of a formative study about how geoscientists collaboratively annotate literature and

Authors' addresses: Shao Zhang, shaozhang@sjtu.edu.cn, Shanghai Jiao Tong University, Shanghai, China; Yuting Jia, hnxxjyt@sjtu.edu.cn, Shanghai Jiao Tong University, Shanghai, China; Hui Xu, xhui_1@sjtu.edu.cn, Shanghai Jiao Tong University, Shanghai, China; Ying Wen, ying.wen@sjtu.edu.cn, Shanghai Jiao Tong University, Shanghai, China; Dakuo Wang, dakuo.wang@ibm.com, IBM Research, Cambridge, Massachusetts, United States; Xinbing Wang, xwang8@sjtu.edu.cn, Shanghai Jiao Tong University, Shanghai, China.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2022 Association for Computing Machinery.

Manuscript submitted to ACM

extract and aggregate data, we proposed DeepShovel, a publicly-available AI-assisted data extraction system to support their needs. DeepShovel leverages the state-of-the-art neural network models to support researcher(s) easily and accurately annotate papers (in the PDF format) and extract data from tables, figures, maps, etc. in a human-AI collaboration manner. A follow-up user evaluation with 14 researchers suggested DeepShovel improved users' efficiency of data extraction for building scientific databases, and encouraged teams to form a larger scale but more tightly-coupled collaboration.

CCS Concepts: • Human-centered computing → Human computer interaction (HCI); • Computing methodologies → Artificial intelligence.

Additional Key Words and Phrases: Human-AI Collaboration, Team Collaboration, Data Extraction, Scientific Literature Processing, Geoscience

ACM Reference Format:

Shao Zhang, Yuting Jia, Hui Xu, Ying Wen, Dakuo Wang, and Xinbing Wang. 2022. DeepShovel: An Online Collaborative Platform for Data Extraction in Geoscience Literature with AI Assistance. 1, 1 (February 2022), 26 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

The shifting of the data-driven research paradigm raises new requirements for researchers to build *scientific databases* [22] in many disciplines, including Geoscience [7], Medicine [6], Biology [3], Chemistry [44]. Scientific database is a collection of structured and verified research results that consists of various numeric, word-oriented, or image-organized data, which plays a central role in data-driven research [14]. In this paper, we focus on the need of constructing scientific databases in geoscience, which can help geoscientists to discover unknown phenomena and novel insights in Earth [17, 18, 54]. To obtain enough high-quality data, geoscientists often review a large amount of published literature [33, 39, 40] (generally are PDF documents), from which they locate and extract useful data (e.g., tables, figures, maps, etc.) to construct the scientific databases. However, it is nearly impossible for a single research team to manually collect and organize the scattered data from hundreds of thousands of documents, and the number is still increasing [43]. Although some larger research teams may have more workforce, without a well-designed computer-supported cooperative work (CSCW) platform, they still need to spend lots of effort to process a sufficient amount of literature and extract enough data for the scientific database. Furthermore, a larger research team may face more team collaboration and coordination challenges such as synchronizing the work process and resolving conflicts. Because of these challenges, constructing a scientific database using the data extracted from a large number of the literature often takes several years with a large workforce. This is a huge obstacle to the advancement of research.

With the development of artificial intelligence (AI), a few research teams have recently begun to explore building scientific databases with the help of AI. Researchers have proposed several fully-automated information extraction systems (e.g., DeepDive [62] and Fonduer [60]) with a goal to efficiently process documents in the PDF format and extract data. Although these end-to-end fully-automated extraction systems may reduce the burden of manual document processing for researchers, they suffered the limitation of insufficient data extraction accuracy [49]. Furthermore, these systems require a huge amount of annotated data before training the AI extraction model, which is a well-established challenge and often infeasible in real-world to collect sufficient training data for multi-context deep-learning models [20]. Unlike many general annotation practices in computer vision and natural language processing [10], the data annotation and extract task in geoscience requires deep domain expertise [26], making it impossible to obtain these labeled data using the general crowd-sourcing platforms [42]. Consequently, researchers would have to pay extra effort for labeling and cleaning the training data to make the AI model work, which may take even more time than they manually extract data without using an AI. We argue that instead of designing a fully-automated end-to-end solution,

a human-AI collaborative and interactive system may be the solution to address these problems. Geoscientists can perform the data extraction activity as they used to, and AI can train itself with these user-labeled data and then make suggestions to the user in the future. Together, the human-AI team can accurately extract data at a much lower cost.

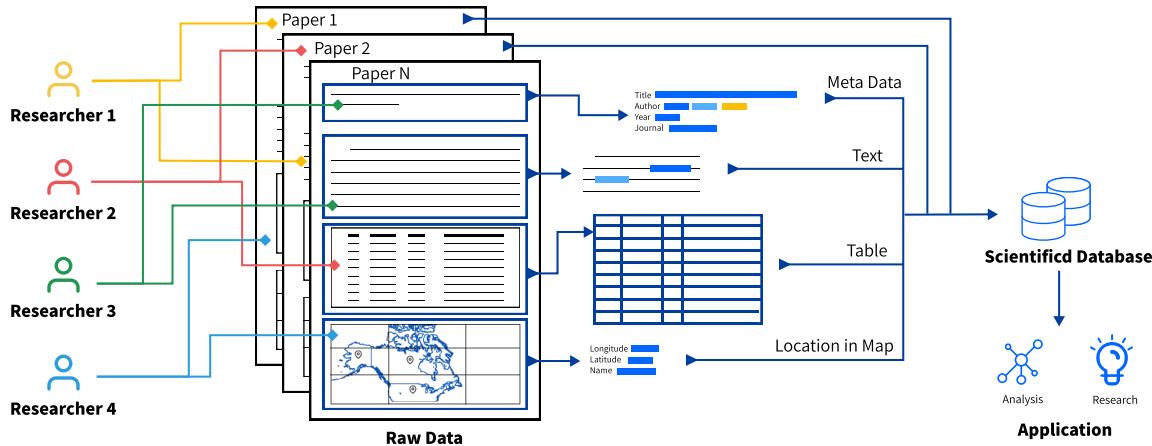


Fig. 2. The collaborative data extraction from scientific literature for big data-driven geoscience research.

To explore how geoscientists extract data from literature and then construct a scientific database, we recruited geoscientists from the Deep-Time Digital Earth (DDE) program [36] and conducted a formative study to understand the problem space and identify user requirements.

The DDE program is a data-driven discovery program in geoscience with a goal to aggregate the geoscience data and to facilitate data-driven discovery for understanding Earth's evolution [57]. We designed and distributed a survey to DDE geoscientists, and then invited some of the respondents to participate in our semi-structured interview sessions. In total, we collected 119 questionnaires, of which 106 are valid, and interviewed 12 researchers with different roles in the team. The formative study reveals the workflow of these geoscientists (Figure 2) and suggests three user needs:

- **N1:** Geoscientists need a more efficient way to collect and extract data from the documents.
- **N2:** Geoscientists need a platform to help them conduct team collaborative data extraction.
- **N3:** Geoscientists need a collaborative platform to organize their group research.

Based on these findings, we then designed DeepShovel, an AI-based collaborative data extraction platform to assist geoscience research teams to complete the work of data extraction, data aggregation, and scientific database construction. DeepShovel provides a user-friendly interface and experience design following the human-AI interaction design guidelines [4], so that users even without any AI backgrounds can also work comfortably with it. DeepShovel has been deployed for one month and there are already 253 users from 36 geoscientist teams within the DDE program using it in a daily basis. More than 240 projects and 46,000 documents have been processed for building scientific databases. We further recruited 14 existing users from 9 teams and conducted a user study to understand their user experience of the system.

This paper makes the following contributions:

- We conducted a needs-finding formative study to explore geoscientists' workflow and challenges in the process of data extraction from literature for building scientific databases, from which we proposed design suggestions for building the AI-based data extraction tools.
- We developed an online team collaborative AI-assisted data extraction platform and have successfully applied it in the process of building scientific databases in geoscience research.
- We proposed a human-AI team collaboration framework making the collaboration process between humans and the AI system more productive and having a better performance, which may generalize beyond the geoscientist data extraction use case.

2 RELATED WORK

2.1 Team Collaboration in Building A Scientific Database

As big data-driven research becomes commonplace, an increasing amount of scientific databases are needed, and building a reusable database becomes a challenge for researchers. In the field of geoscience, many researchers have pointed out that the process of data preparation is vital throughout the research process [49]. Many geoscientists think that building a system for data extraction access is necessary for a scientific database [15].

In the past, many researchers in geoscience tried to build an integrated platform from data collection and storage to data analysis, which promoted the sharing of geoscience databases and big data research. Chronos [9] is a community facility addressing the needs of geoinformatics and providing simultaneous and seamless integration of hosted and federated databases with analytical and visualization tools. Paleostrat [47] is designed as an infrastructure platform for Earth Science researchers and teachers, which serves the community by enhancing the research and education process. However, they spent too much time designing schemas instead of trying to make a user-friendly interface for geoscientists to easily extract desired information, which makes it impossible for researchers to cooperate with their team in the data extraction process. It leads to the conclusion that they do not have enough data to support their operation.

Some researchers noticed and analyzed the problems in team collaboration that existing platforms are facing. Hoeppe [22] mentions that the large scientific databases often require large teams to extract and clean the data, and it is a complex task that needs to be done through CSCW. Schmidt [45] demonstrates that task distribution, allocation, and interrelating of 'distributed individual activities' are some important issues in team collaboration. Steinhardt and Jackson [48] think that these activities require the support of infrastructure. Finholt and Olson [19] concluded three core capabilities of technologies supporting scientific collaboration: linking people with people, linking people with information, and linking people with facilities. However, they only provided some theoretical analysis and solution principles, and no implementation details and practicable platform have been proposed yet.

In this paper, we aim to reveal empirical understandings and build systems to support researchers' data extraction work with their teams. Specifically, we investigated a research 1) to understand how researchers cooperated with their teams extracting data from scientific literature today, 2) to implement a system with AI assistance prototype to support their work, and 3) to explore researchers' feedback and design implications after they try it out.

2.2 Data Extraction in Scientific Literature

With the development of artificial intelligence, more and more researchers pay attention to the data and phenomena in the published literature. Artificial intelligence methods are introduced in this process to help promote the collection of Manuscript submitted to ACM

data from the published literature and build a scientific database in biology [29], chemistry [41, 50], and many other disciplines [8, 23, 56]. In addition, there are also some related works focusing on literature retrieval and dissertation comprehension to extract information from scientific literature [32, 52, 53].

In the field of geoscience, GeoDeepDive [63] is a widely-used toolkit that adopts natural language processing (NLP) technology to process and analyze the literature end-to-end. Peters et al. [38] use GeoDeepDive to study the evolution of stromatolites in shallow marine environments. GeoDeepDive is a case study of DeepDive [62] which stopped update as of 2017 and became a digital library named xDD [46]. Due to the use of supervised learning and end-to-end extraction methods [21], the lack of labeled data has introduced the problem of insufficient data accuracy. Moreover, these methods can only process text [5, 16, 35], use NLP to analyze the tables with the content [21], and do not store the tables as structured data for broader big data analysis.

However, most existing works on data extraction in the scientific literature using artificial intelligence require case customization and do not have a user-friendly user interface for researchers, leading to obstacles for researchers who do not have a artificial intelligence background [28]. Wilkinson et al. [59] demonstrate that scientific data needs to be findable, accessible, interoperable, and reusable. Furthermore, data extracted by the customized case is not reusable, resulting in low data utilization. We believe that retaining structured data can increase data reusability and avoid duplication of data extraction, which is why we consider the idea of Human-AI Collaboration. We study the current workflow in which geoscientists manually extract data from the literature. Based on this workflow, DeepShovel adds the use of AI assistance and human-AI collaboration thinking into the data extraction process. DeepShovel can improve efficiency and data reusability and avoid the lack of precision caused by automatic methods.

2.3 Scientific Document Processing and Human-AI Collaboration

For scientific literature parsing, there are some toolkits and packages like Grobid [1], Science Parse [51] and PDFFigures 2.0 [12]. However, on the one hand, they can only decompose the paper, but still cannot form data that can be used for research. On the other hand, the lack of graphical interfaces means it is difficult for researchers who are only average computer users to use these tools and carry out further data processing easily. Some other interactive document processing tools take ease of use and user-friendliness into account. ABBYY FineReader PDF [2] is a desktop application for processing PDF documents. It provides an OCR tool that makes the content of a PDF document editable. The extraction processes of text, pictures, and tables in the PDF document is mixed so that the user needs to deal with all the information in process, which causes information redundancy and makes it difficult for users to focus on the data they need. TableLab [58], which is an interactive table extraction system, can help people to extract and train a customized model. However, due to the focus model training, TableLab's user interaction process presents a mixture of table structure and content recognition. It means that users need to pay attention to the correctness of both table structure and content recognition at the same time. Such information overload will increase the difficulty of users' decision-making and introduce the problem of model overload [11, 25]. Moreover, these tools are not designed for scientific literature.

DeepShovel provides users with a friendly graphical user interface and a well-design interactive extraction process between humans and AI to prevent information overload and model overload. Users can extract the structured data they need, and the consideration of human-controlled decision making ensures the accuracy of the data. Such a human-AI collaboration framework and functional design with team collaboration can help researchers without relevant technical backgrounds quickly use DeepShovel to extract data from the literature and quickly invest in research with their teams.

3 USER RESEARCH AND REQUIREMENT ANALYSIS

In our project, we follow the design study methodology [37]. As it suggests, we start analyzing the real-world problems of domain experts and working on creating a system to solve these problems. This section describes how we conduct the user research and what we learn from the user research. In the task definition stage, to understand the relevant technical background of the domain experts and the data extraction involved in their research questions, we used questionnaires to investigate 119 related potential users and selected 9 groups of users among them for 60 minutes of in-depth interviews. Then we analyzed and summarized the commonalities of the user groups and their main tasks and difficulties.

3.1 Questionnaire Survey

We choose to use questionnaire surveys [34] to study, which brings better user characteristics understanding, user requirements acquisition, and system design. Specifically, we hope to use the questionnaire to understand the current status of work in the field of earth sciences on the task of constructing subject-specific databases. We divided the questionnaire into three main parts:

- Basic Information
- Task-related Information
- User's understanding of computer technology

Based on the above disassembly of the question, we designed the initial version of the questionnaire and invited potential interviewees to conduct face-to-face interviews to clarify and eliminate the cognitive bias that the interviewee may have. Finally, we determined the topic and the questions of the questionnaire as shown in Table 1.

Since our research is aimed at a specific field, we did not choose the probability-based or random sampling method for survey research. We used non-probability-based methods, invited relevant researchers in the DDE program, took a voluntary form to participate in our questionnaire survey, and combined in-depth interviews to conduct in-depth research. We used an online questionnaire to overcome geographic barriers and covered more researchers in different regions. A total of 119 questionnaires were obtained, of which 106 were valid.

3.2 In-depth Interview

Based on the questionnaire survey, in order to know details of geoscientists' research and team cooperation, we conducted the in-depth interviews as semi-structured interviews [31]. We organized the semi-structured interview with a framework of questions about their research interests, teamwork, and usage of the data they collected to explore their needs during the workflows. We invited 12 users for 30-60 minutes interviews from the questionnaire survey participants who have different roles in their research team. The interviewees' research projects and roles of their research team are shown in Table 2.

3.3 Results

3.3.1 Challenges for Data Extraction. According to the results of the questionnaire, the main challenge is to help geoscientists without backgrounds in computer extract structured data from unstructured data sources, especially PDF documents.

The PDF is the most significant proportion of the document formats that geoscientists pay attention to. Due to the difference in encoding, version, and source of PDF documents, the structure of internally stored digital information is Manuscript submitted to ACM

No.	Questions
Part 1: Basic Information	
Q1	Gender
Q2	Career Position
Q3	Research Field in Geoscience
Part 2: Task- related Information	
Q4	The progress of Data extraction
Q5	Research Team size
Q6	Team composition
Q7	Tools using in data extraction
Q8	Data source used
Q9	Data source format and proportion
Q10	Method to process PDF files
Q11	Distribution and proportion of data in the file
Q12	Number of database fields
Q13	The number of documents needed to build the database
Q14	Personal participation in data collection
Q15	Current Workflow of data extraction
Q16	Time of Single PDF processing
Q17	Estimate the time required for the entire collection
Part 3: Understanding of computer technology	
Q18	Understanding of programming
Q19	Kinds of tasks be accomplished by programming
Q20	Understanding of artificial intelligence
Q21	Data set preparation of artificial intelligence task
Q22	Understanding of data labeling
Q23	Kinds of tasks served by data labeling

Table 1. Questions in the questionnaire.

Group No.	Participant No.	Gender	Research Projects	Roles
G1	P01	Male	Magmatic Migration	PhD Student
G2	P02	Male	Geomagnetism and Geoelectromagnetism	Associate Professor
G2	P03	Male	Geomagnetism and Geoelectromagnetism	PhD Student
G3	P04	Male	Paleoclimatology	Associate Professor
G4	P05	Male	Geochronology and Structural Geology	Professor
G5	P06	Female	Paleontology	Associate Researcher
G6	P07	Male	Structural Geology	PhD Student
G7	P08	Female	Evolutionary Biology and Dinosauria	PhD Student
G7	P09	Male	Evolutionary Biology and Dinosauria	PhD Student
G8	P10	Male	Carbonate Sedimentology	Postdoctoral Researcher
G9	P11	Female	Global Detrital Zircon Database	Full-time Data Entry Clerk
G9	P12	Male	Global Detrital Zircon Database	Full-time Data Entry Clerk

Table 2. Demographics of interviewees.

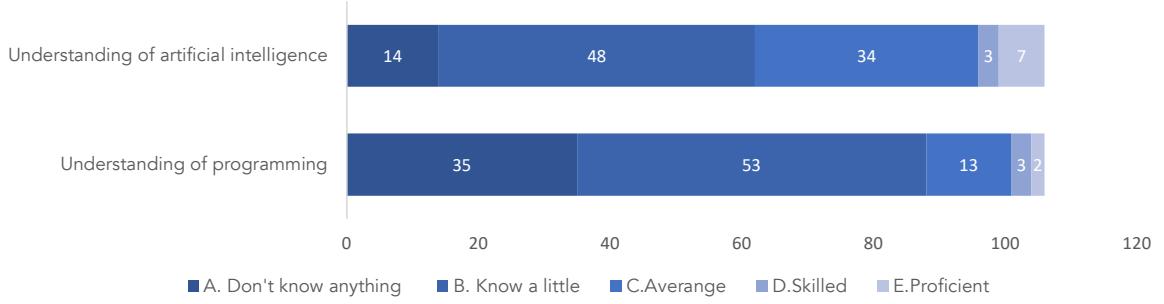


Fig. 3. Results of Q18-Understanding of programming & Q20-Understanding of artificial intelligence.

chaotic and can not be processed automatically by machines. Especially for some PDF files scanned from paper, the scanning quality dramatically affects the results of automated processing. All the interviewees consider such errors unacceptable, as they will directly impact the research results.

In addition, the results of the questionnaires and interviews show that the most valuable data are located in the tables in the article. However, tables in PDF are usually not reproduced very well by ordinary PDF editors, which makes data collection very difficult.

Another problem introduced by completely manual copy-pasting of data is the proofreading of the data. The interviewees P01, P11, and P12 say that they need to spend much time proofreading the data and matching them with the meta information of the corresponding source article in the database.

3.3.2 "We don't know much about programming for data extraction, and we can't utilize state-of-the-art AI models without graphical user interfaces". From the survey research and in-depth interview, we conclude the typical profile of our users that they are average computer users with no / low programming skills. Only P07 is a proficient programmer and can use Python scripts to process the data, but he still mentioned that "Being unfamiliar with artificial intelligence makes me unable to deal with the data in pictures and charts in the papers efficiently". We can find that most participants of the survey also lack understanding of programming/artificial intelligence (Figure 3). Furthermore, some interviewees also report that they always encounter difficulties in data extraction and processing due to the lack of programming skills. We can learn from the questionnaire results that our users have to use a combination of a series of tools to accomplish a task (as shown in Figure 4). They usually use some OCR tools with graphical user interfaces to process PDF documents to make them editable and then manually copy-paste the data they find into Excel. We found that an online collaborative application with a user-friendly graphical user interface is critical for geoscience research teams.

3.3.3 Difficulties of Team Collaboration. Interviewees mentioned that they worked as a team to accomplish data extraction tasks and indicated some problems in team collaboration remain unsolved.

The first difficulty is the storage of data sources related to the distribution of team tasks and the advancement of tasks. We notice in the questionnaire results that respondents generally believe that more than 1,000 articles are needed to construct a scientific dataset. Moreover, we learn from the interviewees that the current methods of researchers sharing literature within teams are still relatively primitive (e.g., copying files and excel sheet records). Using these methods is quite laborious when sharing a large number of files and can cause huge risks of data errors and version conflicts. We believe that building a scientific database is a close-cooperation work but does not have a CSCW system to support it.

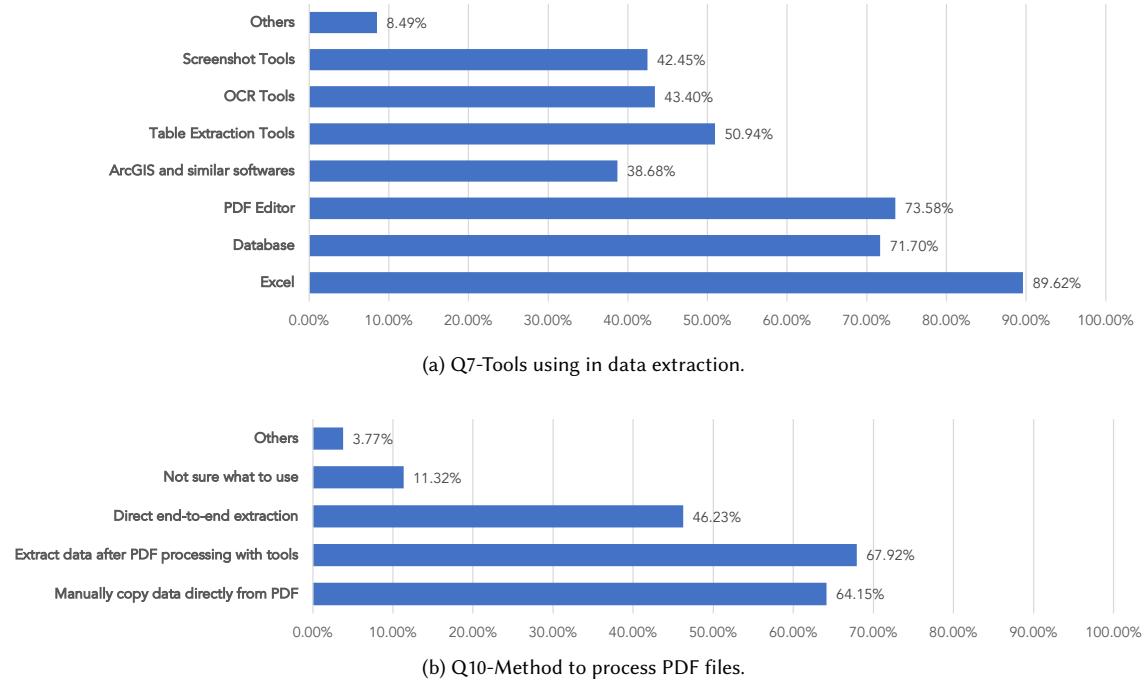


Fig. 4. Results of Q7 & Q10.

Another finding is that Excel is an essential tool throughout the process. Once the data is extracted from the PDF, it will be stored in Excel as we mentioned in §3.3.2. P01 and P13 mention that each team member stores data from his documents into a local Excel file and then merges it with other people's Excel data files in the team. There are great similarities in the ways of collaboration among different teams, which are all primitive and loose.

3.3.4 Summary of the Workflow. Although researchers from different sub-fields of geoscience may study different research questions and focus on different data, their research processes and workflows are basically the same. In this paper, we focus on the data extraction process of their research, as shown in Figure 5.

Based on the results of interviews and questionnaires, we described the data extraction process and defined the tasks in the process. These tasks are grouped and detailed by the main workflow steps in the following list:

- **T1-Problem Definition:** Define the research problem and the structure of the scientific database that needs to be built,
- **T2-Search:** Search for the paper that may contain the data about the research problem,
- **T3-Browser:** Quickly browse the article to find data that is needed,
- **T4-Meta Information Extraction:** Record the literature's meta information for tracking the data,
- **T5-Detail Data Extraction:** Extract data from different parts of the literature,
 - **Data extraction from table:** Get the data in the table and fill in the Excel file prepared in advance cell by cell,

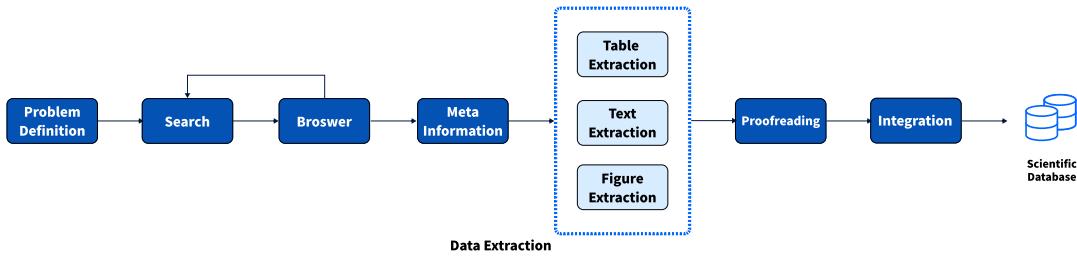


Fig. 5. The team's workflow of collecting and building a database from geoscience literature.

Role	Responsibility	Possible person(s) for this position	Participate in the task flow
Owner	Determine the research topic, Collect relevant literature, Determine database structure, Check the quality of the data.	Research team leader	T1 & T2
Manager	Distribute tasks, Manage task progress, Proofread the extracted data, Integrate data.	Senior researchers or Senior members of the team	T3, T6 & T7
Member	Specific data extraction.	Junior researcher or Professional trained data entry clerk	T4, T5 & T6

Table 3. Teams structure and related tasks for each role.

- **Information extraction from text:** Search the full text with keywords to locate the data, fill in the data in the Excel file after finding it, and repeat until all the data is found,
- **Information extraction from figure:** Restore the corresponding information from the charts, such as obtaining the latitude and longitude of a marked point from the map,
- **T6-Proofreading:** Check and proofread the data to ensure the data is accurate (usually done by the member and the team manager together),
- **T7-Data Integration:** Integrate the data extracted from each paper (usually stored in a bunch of Excel files) into the final dataset.

We also summarized the usual structure of teams building scientific databases and the related task for each role as shown in Table 4. We also want to mention that some interviewees said they built the scientific database individually because not much literature needed to be processed. Meanwhile, in some small teams, the responsibility of the role of "Member" might be taken by the role of "Manager" due to the lack of manpower.

3.3.5 Design Requirements. In this workflow, we found that there are three levels of requirements influencing the efficiency and experience:

- **R1:** Quickly and accurately extract data from PDF files and form structured data,
- **R2:** Help proofread and integrate the data extracted by each person to build a database in multi-person collaborative extraction,
- **R3:** Share tasks' resources (raw data) and tasks' progress across the team.

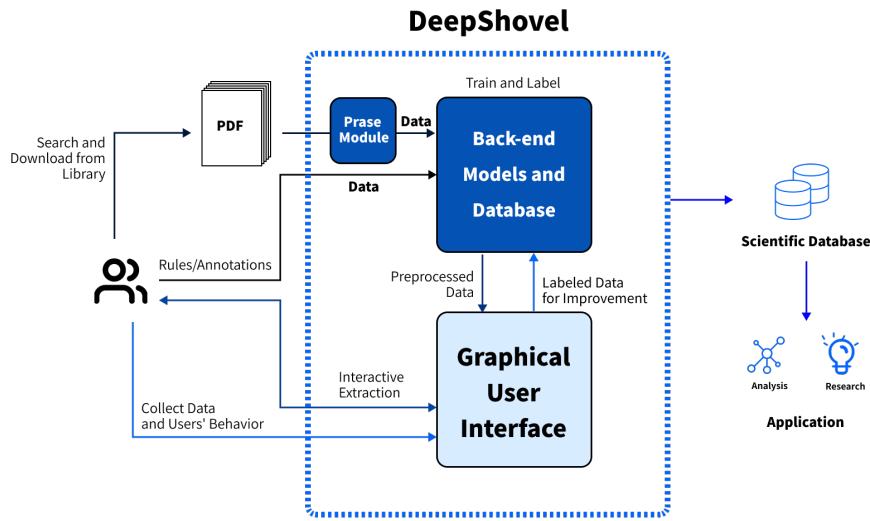


Fig. 6. System overview.

4 SYSTEM DESIGN

Based on user research and requirement analysis, we solve these problems using a human-AI collaboration system design to extract data easily from PDF and have a better team cooperative experience. For each task mentioned in §3.3.4, we design and implement some artificial intelligence modules. According to human-AI collaboration thinking, the system can collect the data for artificial intelligence model training while assisting humans in finishing the tasks. Such collaboration motivates people to participate and allows the machine to obtain enough information to improve. The models implementation details are shown in §4.5.3.

4.1 System Overview

As shown in Figure 6, DeepShovel consists of: (1) a interactive graphical user interface (see Figure 1) including data extraction, document management, team management and data integration (D in Figure 1); (2) a back-end parse module to pre-process the PDF format files; and (3) some back-end artificial intelligence models supporting data extraction and integration functions.

4.2 Human-AI Collaboration for Data Extraction

When users open a file from Project File List to start their work, they will enter the data extraction interface (Figure 7). In the data extraction interface, users can switch the different tabs (e.g., Meta, Text, Table, and Map) in the area F1. The details of each function are in the following sections.

4.2.1 Meta Information Extraction. For the task **Meta Information Extraction**, we develop a module to automatically extract the title, author(s), journal/conference, and other meta information from the PDF file. As shown in Figure 7, users can edit and save the meta information that can be joined to the output dataset (refer to §4.3).

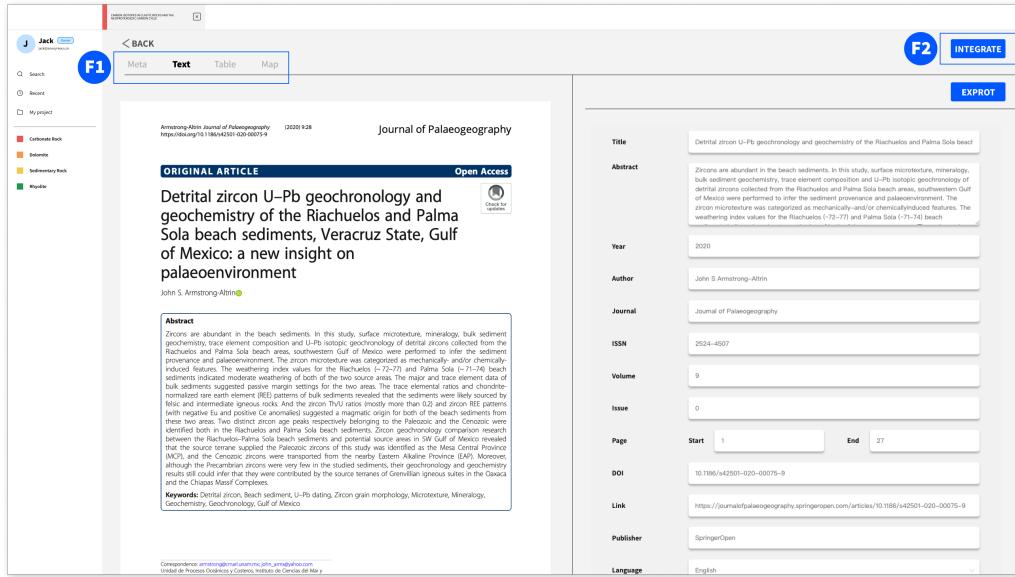


Fig. 7. UI of meta information extraction.

4.2.2 Text Extraction and Annotation. We use weak-supervision learning models and rules to help highlight the focused keywords **in texts** and the samples' features to help them add these words into the database. For example, we have a dictionary of eras' names to highlight the era mentioned in the PDF. Users can also annotate a keyword via mouse selection when they switch to the edit mode (F3) and select a label (F4) as shown in Figure 8. The keywords can be added to the output database (refer to §4.3). Users can choose to show or hide some labels, which are set at the project level as shown in Figure 12b.

4.2.3 Table Extraction. To help user extract the data in the table, we develop the **Table Extraction** function (Figure 9). In this part, we separate the task into three steps: 1) Locate the Table with the assistance of AI; 2) AI recognizes the table's structure, and users assist for a better result; 3) AI recognizes the table's content, and users can edit for final accuracy. In each step, the artificial intelligence models we design will help people to easily get the result and collect the users' adjustments for model training(see Figure 16). From the user's perceptive, the first step is adjusting where the table is in the F5 area or drawing a new area as a table, then starting the recognition of structure. The next step is to adjust the structure that the system advised (F6). The system provides 'add and delete column/row' and 'merge or split cell function'. After structure recognition, users can start the content recognition and edit the content in each cell (F7).

4.2.4 Map Recognition and Location Extraction. For collecting the location of a sample, we provide a module that can **recognize maps** and calculate the latitude and longitude of each point on the map (Figure 10). Users can draw an area (F8) that contains the map and mark a point by right click (F9). The latitude and longitude will automatically be saved in the table (F10) as shown in Figure 10 and can be joined to the output dataset (refer to §4.3).

This screenshot shows the DeepShovel interface for extracting text from a scientific publication. The main area displays the abstract and introduction sections of an article about zircon in beach sediments. A sidebar on the left lists project members (Jack, Bevert, Myself) and categories (Carbocate Rock, Facies, Sedimentary Rock, Facies). A legend at the bottom defines symbols for various geological features. The right side includes buttons for 'INTEGRATE' and 'EXPORT'. A blue circle labeled 'F3' highlights the text extraction process.

Fig. 8. UI of text extraction.

This screenshot shows the DeepShovel interface for extracting tables from the same scientific article. It displays two tables: 'Table 1: Characteristics of mechanical and chemical features identified on the zircon grain surface in the Río' and 'Table 2: Characteristics of mechanical and chemical features identified on the zircon grain surface in the Río'. The interface includes a 'Target Recognition' section, a 'Structure' section, and a 'Content' section where tables are identified. A blue circle labeled 'F5' highlights the table extraction process. A legend at the bottom defines symbols for various geological features.

Fig. 9. UI of table extraction.

4.3 AI-Assisted Team Collaboration

After the data is extracted step by step, it needs to be integrated into a table to establish a database finally. It involves how the data extracted by everyone in the team can be put into a summary table faster. We designed a single file

Manuscript submitted to ACM

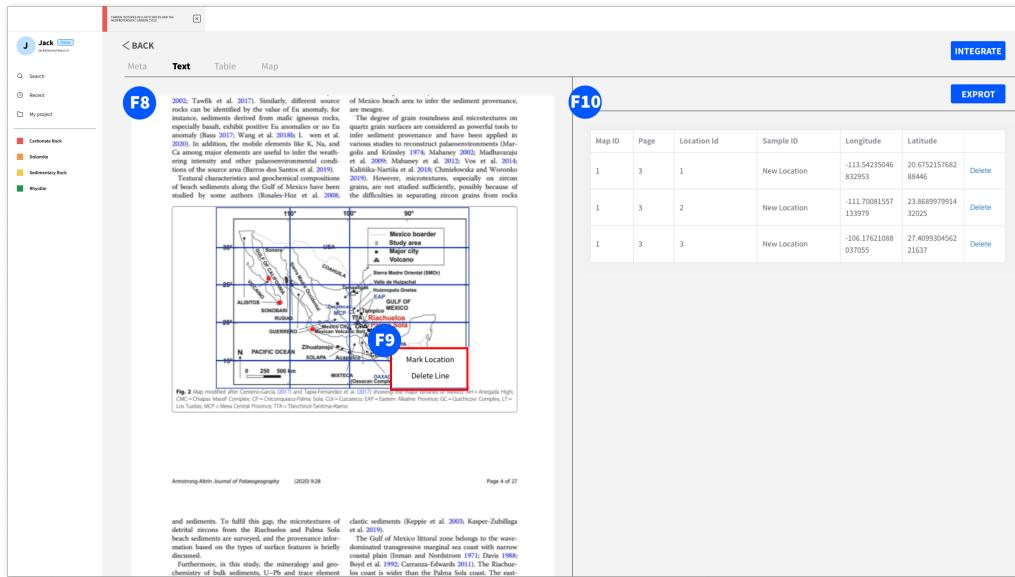


Fig. 10. UI of map recognition and location extraction.

integration and project-level integration with the assistance of AI to adapt to different teamwork modes. The user needs to set the header of the total table on the project page (Figure 12a). Then in the data extraction interface (Figure 7), when users click the Integrate button (F2), the back-end model will process the data in each part, including meta information, tables, location in maps and texts. The result are shown at the F11 area in Figure 11.

After all the data in a single file has been integrated into a file-level summary table, the user can integrate the summary table of each file into a project-level summary table at the Files List interface (F12 in Figure 14), and the result will be automatically downloaded.

4.4 User-Document Management Design for Team Collaboration

As we mentioned in §3.3.3, users need to process a large number of files and need to distribute files to team members for data extraction. Therefore, we provide a user-document management with graphical user interface that can be easily used.

4.4.1 Project and File Management. In order to realize the management of projects, we designed the projects list interface to display the relevant information of each project, as shown in Figure 13. Each project has a file list (see Figure 14), which will show who uploaded the file, who was the last editor, the upload time and last edit time, and whether the file has a principal (refer to §4.4.3). Users can change the project settings, including the text labels, the export dataset headers, and the project description. Considering that the dataset may contain several headers, we provide batch edit for convenience.

4.4.2 Roles Settings and Cross-team Cooperation. To help researchers manage their team in data extraction work, we preset three team roles in the system design: Owner, Manager, Member. The user permissions of the three roles are shown in Table 4.

Data Table:

File ID	Sample ID	Age	Sm	Nd	147Sm/144Nd	143Nd/144Nd	εNd(t)	TDM2
TUrof02H4464652n02n121M0t0	EMR-1C	1.045	1.204					
TUrof02H4464652n02n121M0t0	EMR-1C_1	1.143	1.075					
TUrof02H4464652n02n121M0t0	EMR-1C_2	1.073	1.170					
TUrof02H4464652n02n121M0t0	EMR-1C_3	1.294	1.193					
TUrof02H4464652n02n121M0t0	EMR-1C_4	1.131	0.958					
TUrof02H4464652n02n121M0t0	EMR-1C_5	1.124	1.080					
TUrof02H4464652n02n121M0t0	EMR-1A_3	1.714	1.527					
TUrof02H4464652n02n121M0t0	EMR-1A_4	1.427	1.364					
TUrof02H4464652n02n121M0t0	EMR-1A_5	1.375	1.343					
TUrof02H4464652n02n121M0t0	EMR-1A_6	1.797	1.651					

Table 1: Microtexture and mechanical features identified on the zircon grain surfaces in the Ratchaburi and Palma Sola samples.

Microtexture	Zircon grain Ratchaburi	Palma Sola	Palaeoenvironment ^a
Mechanically-induced feature			
Abraded edge (are)	X	XXX	Aeolian, saltation, collision
Dual oriented zircon (diz)	X	XX	Saltation, collision, short transport
Squeezed zircon with one side broken edge (sze)	XX	X	Aeolian, saltation, collision, short transport, storm record
Cross-cutting fracture (cfr)			
Wavy shear zone (wsz)	X	X	New shear, wave action
Anti-shape shear zone and Lateral step (as)	X	X	High-energy collision, aeolian, glacial zone
Abraided edge (abe)	XX	X	Aeolian, saltation, collision, glacial zone
Rebewelded crosscut fracture (rcf)	X	XX	High-energy collision, aeolian, fluvial zone, nearshore subaqueous
Collision fracture (cf)			
Abraided particle (ap)	X	X	High-energy collision, aeolian, fluvial zone
V-shaped penetration crack (vpc)	X	X	High-energy collision, gouging, fluvial zone, deltaic, subaqueous, surf zone
Straight groove (sg)	X		Littoral zone, wave action, saltation
Chemically-induced feature			
Diagenetic environment feature (dfe)	XX	XX	Diagenetic environment, high in concentration sea water (alkaline fluid)
Circular solution pit (csp)	XX	X	Intertidal zone, diagenetic, percolation of sea water
Gran cavities (gc)			Diagenetic, percolation of sea water
Abraided particle (ap)	X	X	Diagenetic environment, nearshore
Silica particle (sil)	X	X	Starting stage of in-situ diagenesis, nearshore
Abraided particle appears to be silica particle (asp)	X	X	In-situ diagenesis, silica saturated, low energy
Silica flower off and crystal overgrowth (sc)	X	X	Advanced stage of diagenetic environment, silica overaturated nearshore
Abraided particle (ap) ^b	XXX	XXX	Diagenetic, littoral, low energy

^a See notes in Table 1 for details. ^b Note: Abraided particle (ap) and Abraided particle (ap)^b have identical meanings.

Fig. 11. UI of data integration.

Project Settings

Text Table Fusion

- label
- era
- formation
- region
- subject
- thickness
- stratum
- age
- location

Project Settings

Text Table Fusion

Keywords Group	Description
SrNd	Strontrium neodymium isotope
nju	detrital component
e-IODP	International Ocean Discovery Program
mineral-resources	knowledge graph and intelligent prediction c
Magmatite	Magmatite
nannofossil	nannofossil
Early Palaeozoic Marine Biodiversity	Collecting global stratigraphical data
ELIP-UYCB	upper Yangtze cratonic basin S2S system

(a) The settings of text extraction.

(b) The settings of data integration.

Fig. 12. The project settings.

According to the users' description of their current team structure and cooperation in the formative study, the users intend to ensure the original data's controllability and distinguish different research projects (the same team may carry out multiple projects). Therefore, in team management, we ban the modification of the project by the Member role to prevent the original data from being modified. At the same time, to ensure the rigor of the output dataset, we only open the modification of project settings to Manager and Owner.

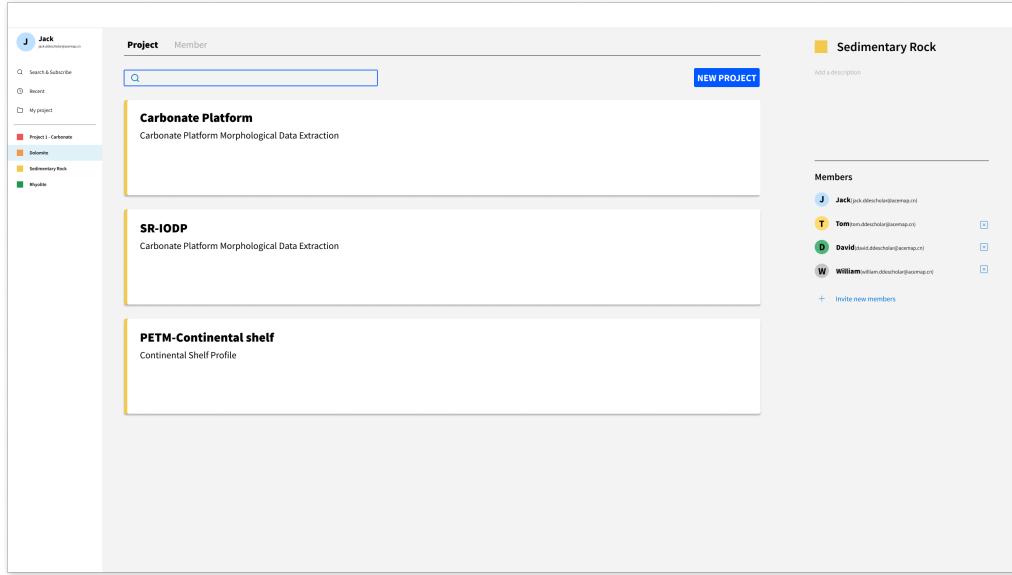


Fig. 13. UI of project list.

No.	File Name	Principal	Latest Editor	Latest Update	Import User	Import Time	Meta	Figure	Table	Text
1	Spectrographic Determination of Ca and Mg in Limestone and Dolomite; GEOLOGICAL NOTES	Take charge	Jack	2021.07.27 14:55	luonu123	2021.06.11 14:25	x	x		
2	Diagnoses of a mixed dolomite-expansive deposit of the Middle Mississippian Sheldene Project, the Canadian Cordillera Range, NE S. China	Tom	Tom	2021.07.28 11:04	luonu123	2021.06.11 14:21	x	x		
3	Porosity Styles of the Middle Field in Western Basin of Southeastern Saskatchewan; ABS TRACT	Tom	Tom	2021.07.02 22:29	luonu123	2021.06.11 14:20	x	x		x
4	Stable isotope and elemental characteristic of sedimentary rocks from a Karstic Krabi River (south east Thailand); useful environmental prospectives?	David	David	2021.06.29 13:26	luonu123	2021.06.11 14:19	x	x		x
5	Cause of dolomization by sheetfold vs. debris flow processes on the adjoining alluvial fans, Death Valley, California	Tom	Tom	2021.06.37 20:15	luonu123	2021.06.11 14:17	x	x		x
6	Paleozoic Rocks of North-Central Nevada I	-	-	2021.06.11 14:16	luonu123	2021.06.11 14:16	x	x	x	x
7	Dolomitic Mountains and the origin of the dolomitic rock of which mainly consist: historical developments and new perspectives	Take charge	William	2021.07.27 17:17	luonu123	2021.06.11 14:15	x	x		
8	Lithologies and cyclicity of dolomites evaporite facies in the rift margin of the Levant in the Late Tertiary, Northern Jordan, South Israel	Take charge	Tom	2022.01.06 22:35	luonu123	2021.06.11 14:13	x	x		
9	Genesis of Mississippi Valley Type Lead and Zinc Ore Deposits and Consequent Exploration in China and Australia	Take charge	Jack	2021.06.29 14:19	luonu123	2021.06.11 14:12	x	x		x
10	Origin and properties of the Late Precambrian Gneissitic Dolomites of Pithoragarh, Kumaon Himalaya, India	Take charge	David	2021.06.25 14:13	luonu123	2021.06.11 14:11	x	x		x

Fig. 14. UI of files list.

Our cross-team collaboration design is based on research in the CSCW community on patterns of information exchange among teams in research work [55], and on respondents' descriptions of this behavior. Information exchange Manuscript submitted to ACM

Role	Add/Remove Manager	Add/Remove Member	Add/Delete Project	Import File	Project Settings
Owner	√	√	√	√	√
Manager	-	√	√	√	√
Member	-	-	-	√	-

Table 4. The user permissions of each role.

is required among teams due to knowledge sharing. In order to meet the common scenario of cross-team collaboration in research, we allow users to join different teams. Team member removal does not affect his/her past actions.

4.4.3 Lock and Principal Mechanism. Online collaboration systems often encounter the situation of simultaneous editing. Considering the particularity of system functions and the interaction process with the back-end model, we designed a file locking mechanism to prevent more than one user from operating on the same file. Based on file lock, we implement a principal mechanism. Users can click the "Take Charge" button in the list to choose to be the "principal" in charge of a file. Then the file can only be operated by the principal user, and other users can only read it. The user can release the file permission at any time.

4.5 System Implementation Details

4.5.1 Web Application. The front-end interactive single-page web application of DeepShovel is developed in Vue.js and hosted with Nginx. The web-based design of DeepShovel gives it the ability to run in the web browsers on a variety of platforms including desktops, laptops, tablets, and smartphones. The use of Vue.js and the design of single-page-application bring extreme load speed similar to native apps and consistent user experience across devices and platforms.

The back-end API service of DeepShovel is implemented with Python and FastAPI framework. The asynchronous coding design makes it possible to achieve higher concurrency with a minimal resource occupation so that it can support more users at the same time. We adopt a master-slave backup MySQL database to store documents and extracted data, which ensures data security and efficient reading and writing. In terms of user system security, we only store and bcrypt hashed passwords to ensure that users' plaintext passwords will not be stored and leaked. Moreover, the HTTPS protocol is applied to the whole system of DeepShovel to ensure the security in network communication.

4.5.2 Document Management & Retrieval. All the literature uploaded into DeepShovel are all automatically parsed with Grobid [1] and Science Parse [51]. The meta information of papers (e.g., Title, Author List, Abstract, Venue, and Year) is extracted and indexed with Elasticsearch. Then all the fields could be utilized for searching and retrieving the documents. Moreover, to better browse and manage literature, the document list could be filtered with the principal user as well as the import user and sorted by title, import time, and latest update time. To go further, each user could get "My File List" containing only documents taken charge by him and "Recent File List" containing his recent viewed documents, which allow users to obtain the documents most important to them and simply continue their respective workflows.

4.5.3 Data Extraction. Generally DeepShovel extract data from these parts from literature:

- **Meta Information Extraction** For each uploaded document, DeepShovel uses multiple parsing tools (e.g., Grobid, Science Parse, and PdfFigures 2.0) to independently extract its meta information and mix all the information with a voting mechanism.
- **Table Extraction:** First, DeepShovel uses an object detection model Detectron2 [61] trained on TableBank [30], a benchmark dataset for table detection, to detect the region of tables. Then for each table, a series of rules are adopted to locate each cell within it. Once users confirm the cell structure of a table, Tesseract [27] will be applied to detect the text in each cell and establish the final digitalized table.
- **Text Extraction:** To extract academic entities from papers with the format of PDF, DeepShovel first utilizes PDFFigures 2.0 [12] to parse each text section from the original files. Then some rules and the natural language processing library spaCy [24] are adopted to automatically extract entities of different types from the parsed text sections.
- **Map Recognition and Location Extraction:** Users can box the region of any map they care about. Then DeepShovel will detect the longitude and latitude labeled at the margin of the map and determine the coordinate range of the entire map. Then if users click any location on the map, the exact coordinates of the location will be automatically calculated and recorded.

5 EVALUATION

We conducted a user study to evaluate DeepShovel. The study examined the following research questions:

- **Q1: Can DeepShovel cover all the tasks of the data collection from literature?**
- **Q2: Can researchers successfully obtain the data that can build the database?**
- **Q3: Can DeepShovel help researchers better collaborate in teams?**

5.1 Participants

We invited 14 users (U1-U14) from 9 different teams (T1-T9) currently using DeepShovel to join our evaluation. They started using DeepShovel at the same time (according to the registration time). The demographic characteristics and backgrounds of the participants are reported in Table 5.

5.2 Procedure

Each user study session lasted around 45 minutes and was conducted remotely via Tencent Meeting due to the COVID-19 pandemic. Participants accessed DeepShovel using the browser on their computers and shared their screens with the experimenter. All sessions are recorded on video.

Considering the participants are all the current users of DeepShovel, the experimenters only briefly introduced the DeepShovel and the study. In order to evaluate the system coverage of extraction tasks, we first asked the participants to extract the data of a representative article selected on their own about their research with DeepShovel. As mentioned in §4.2, we provided meta information, table, text, and map extractions. We asked participants to process the file with their usual workflow and observed whether their operating procedures aligned with our system design assumptions. After the data extraction part, each participant filled out a post-study questionnaire. Then we had a 15-minute semi-structured interview with each group of participants about their team collaboration experience with DeepShovel.

TID	UID	Gender	Role in Team
T1	U01	Male	PhD Student
T1	U02	Female	PhD Student
T1	U03	Male	PhD Student
T2	U04	Female	PhD Student
T2	U05	Male	Associate Professor
T2	U06	Male	PhD Student
T3	U07	Male	Postdoctoral Researcher
T4	U08	Female	PhD Student
T5	U09	Female	PhD Student
T6	U10	Male	PhD Student
T6	U11	Female	PhD Student
T7	U12	Male	PhD Student
T8	U13	Female	PhD Student
T9	U14	Female	Associate Researcher

Table 5. Demographics of user study participants.

5.3 Results

All 9 groups of participants have shown us how they use DeepShovel to extract data from the scientific literature. The time used on each participant is about 10 minutes. There is no apparent deviation between the operation process of each participant and the process we envisioned. However, the difference among their research interests means that the data they need to extract is distributed differently in the article, leading to the difference in the function order. *Meta Information Extraction* is a function that every participant use. Among other functions, most users will use *Table Extraction* first and then *Text Extraction*, while some users will only use *Text Extraction*. A small number of users will use the function *Map Extraction*. Since teamwork is a long process and it is difficult to observe directly, we conducted an assessment in the questionnaire and learned the specific situation in the interview. We will discuss the results from the post-study questionnaire and the interview in the rest of the section.

5.3.1 Post-study Questionnaire. We asked all participants to fill out the post-study questionnaire to rate the usability, usefulness, and user experience of teamwork. The questionnaire uses a 10-point Likert scale from "strongly disagree" to "strongly agree". Note that two participants failed to fill out the post-study questionnaire. The results of the questionnaire are summarized in Figure 15. Specifically, DeepShovel scored on average 6.75 (SD=1.42) on "I am satisfied with DeepShovel", 6.33 (SD=2.83) on "DeepShovel is easy-to-use", 7.91 (SD=1.44) on "I'd like to continually use DeepShovel in the future", 8.00 (SD=1.91) on "I am willing to recommend DeepShovel to other researchers". And about team collaboration, DeepShovel scored on average 6.67 (SD=1.75) on "DeepShovel improved the efficiency of my teamwork", 7.3 (SD=1.59) on "DeepShovel improved the efficiency of my data integration". We find that the score of satisfaction is lower than the score of willingness to continue to use, and the overall score is quite different. We believe that this is due to different data extraction tasks caused by different research projects (see Table 5), and we will focus on explaining this problem in the interview results.

5.3.2 User Experiences and Feedback. We have a 15-minute semi-structured interview with each group of participants at the end of our user study. The post-study questionnaire results are discussed, and the questions about user experience (including team collaboration experience) are asked in the interview, which brings some important findings:

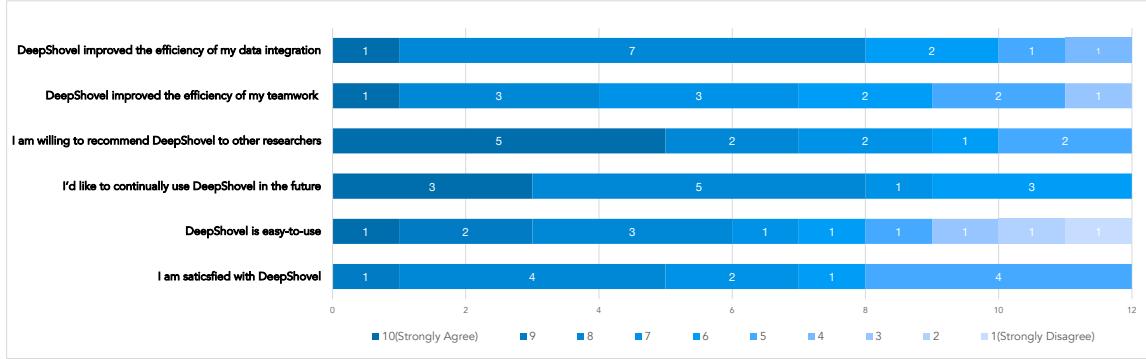


Fig. 15. Results of the post-study questionnaire.

DeepShovel brings the possibility for geoscientists to form a team to build the database collaboratively. Actually, some participants in the user study did not form a team to carry out this work before because of the difficulty of data extraction and the lack of effective team collaboration tools and workflows. U08 from T4 said "My supervisor only considered forming a data extraction team to advance big data-related research after learning about DeepShovel". U14 from T9, U12 from T7, and U09 from T5 also mentioned a same situation in their research. We think this shows that DeepShovel has helped researchers to promote data-driven research to a certain extent, prompting them to conduct deep-time and deep-space research. It also proved that our UI design and AI assistance for teamwork are effective. On the other hand, there are cases of non-team use, such as U07, U04, and U05 from T2 (although they are from the same team, their work does not overlap, and they do not have teamwork tasks). Even though they do not need the functions for team collaboration, they are also satisfied with the system because it effectively reduces the complexity of data extraction. No matter it was for team use or personal use, the participants gave a relatively satisfactory score for the AI Assisted Team Collaboration, which is the design that assists in the integration of data. For example, U01 from T1 said "The integration of data from a single table into a comprehensive table allows me to enter all the data into the database faster".

Different data requirements lead to different user experience. The reason for users' frustration with lower satisfaction in the questionnaire is some technical limitations, such as the parsing problem introduced by PDF files and the accuracy of text extraction. For example, U07 from T3 mentioned "The data I care about are always in the text, but the region name recognition is very bad". U14 from T9 said "The inability to customize tags prevents me from labeling the words I care about". Participants were generally satisfied with the table extraction module. When many participants mentioned the table extraction, they all said "it has greatly improved my work efficiency". U13 from T8 said "In the past, I thought that completing the data collection work was a distant matter, but now DeepShovel makes me feel that it is possible and close at hand". Participants also mentioned that the level of use of different modules affected their satisfaction scores to a certain extent. When participants mainly use table extraction to complete tasks, in other words, when the data they care about is mainly in tables, they are more satisfied with DeepShovel.

Different teams have different collaboration ways when using DeepShovel. U01 from T1 said that they have 18 people in the team. They fully used the team management system and basically aligned with our vision. However, the situations are mainly different when they are in a smaller team. U10 from T6 said "We did not have distinguished team roles, but collected documents together and divided the tasks equally". Therefore, they creatively use the take charge function to

mark whether the task is complete instead of dispatching the task. We think this shows that the collaborative design of DeepShovel can adapt to the use of different team sizes, while also applies to the individual use scenarios as mentioned before. Furthermore, some functions should be further expanded to mark the completion of the task, which is a function that we do not currently provide.

Cross-team cooperation took place in the U13 and U14 teams. In order to distinguish the groups of participants, we did not indicate this phenomenon in Table 5. U13 mentioned that she participated in the research work of three real-life teams, which is common in scientific research. However, she also mentioned that some of her personal work is also stored in DeepShovel for convenience. She said "This actually caused confusion in file management to some extent". We think this is also an interesting phenomenon that we never imagined. In the future, we should consider improving the convenience of cross-team collaboration to avoid this kind of file management confusion.

6 DISCUSSION

The results from our user study and the practical application suggest that geoscientists can successfully collaborate with DeepShovel to extract data from scientific literature and integrate the data into databases. DeepShovel also performed well in helping researchers' cooperation in their team. In this section, we discuss the lessons we learned and the design implications of our work.

6.1 Data Extraction from Scientific Literature with Human-AI Collaboration

6.1.1 Why End-to-End is Not a Good Choice? Due to the accumulation of errors in end-to-end approaches, the final outcome might be unacceptable. Meanwhile, manually checking and correcting these errors is a very tedious and difficult job. We would like to claim that we do not believe today's AI technology can build a "fully automated" system to replace researchers in data extraction. To ensure the quality of the database used in further research, researchers still have to clean and correct data manually. We think that building a human-AI collaboration solution with the appropriate level of automation would be a better way to solve the problem so that the human and AI can jointly iterate, improve and complete the data extraction. The user makes the final decision of all data extraction, and AI fully follows the user's instructions in this interaction process to ensure the accuracy of the data. For example, in the table extraction process, AI only recommends where the tables are, and users decide which table they would like to extract.

6.1.2 Decision Making in Human-AI Collaboration. In DeepShovel, AI only suggests the tables' structure and content, and users can edit this information and decide on the final output. The user will take turns interacting with different models during a table extraction, each model explaining its own task to the user. Such an interactive form can first effectively help users understand the role of the model, allowing users to make decisions more intuitively as to whether they need to be modified. Secondly, due to the step-by-step interaction (Figure 16), information overload is avoided, and the user can focus more on the current decision. Considering the interaction process between models and humans, preventing information overload and model overload, and making the interaction process as disassembled as possible to fit into the human decision-making process are two critical issues to be considered.

6.2 Team Collaboration and Task Distribution

6.2.1 Fine-grained Functional Division. There may be many different forms of task assignment in multi-player cooperative tasks. If the framework design restricts task division, it will cause damage to the original team operation mode. When each user performs a data extraction task in a single PDF document, DeepShovel follows the rules of

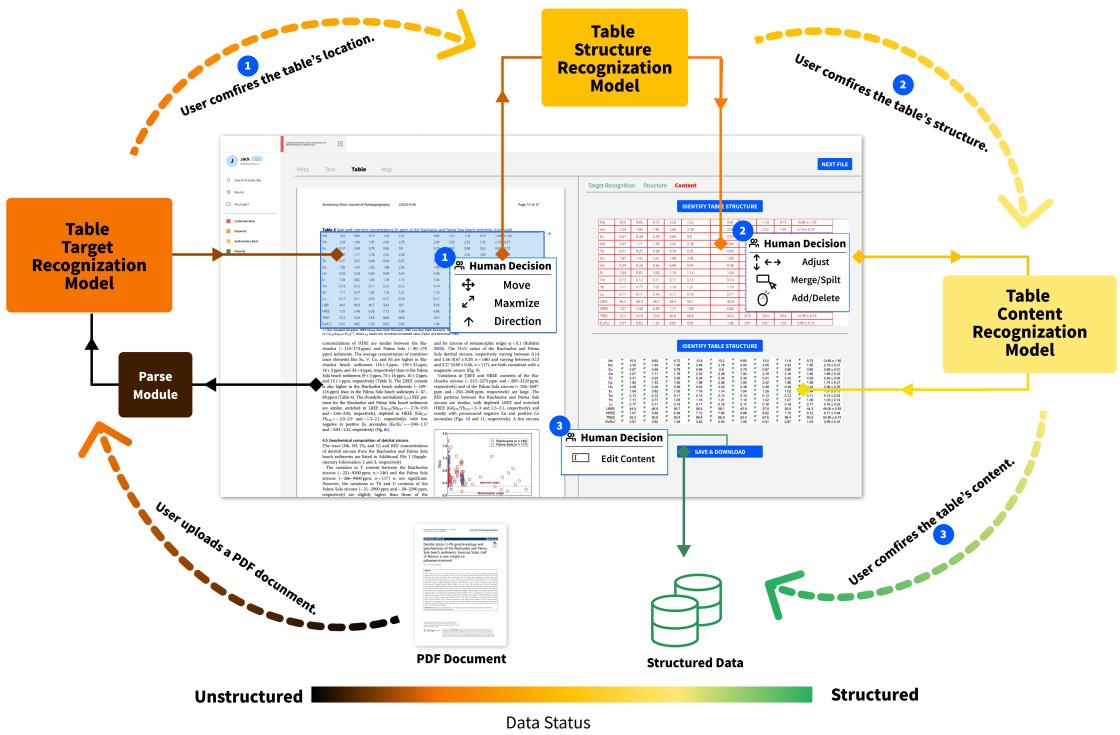


Fig. 16. The step-by-step interaction in table extraction.

the distribution of target data in the document and breaks down the task into meta information, text, table, and map extraction, and each subtask is served by a different model. The models are introduced by task division so that users do not interact with multiple models at once. Such a design can avoid users from processing information from three different modalities: image, text, and table at one time. This splitting of functional modules based on data distribution brings greater flexibility to team collaboration while reducing the pressure on online services. In DeepShovel, the team can divide labor according to projects, documents, or even tasks, which adapts to the existing team cooperation methods to the greatest extent. We believe that when designing related collaboration platforms in the future, different forms of team collaboration should be considered in the process design.

6.2.2 AI Assisted Team Collaboration. Our research found that no tool exists in current team data integration, and geoscientists have to integrate data using Excel manually. We think this is an essential step after data extraction, but there is still no very reasonable solution for some researchers who are just average computer users. However, there are still problems with this design. Some users report that data cannot be merged correctly because different authors have different descriptions of some fields when drawing tables, and the information in the text cannot be merged.

6.2.3 Information Sharing in In-team Collaboration. In team collaboration tasks, information sharing within the team is crucial to the efficiency and quality of task completion. We studied the teamwork of researchers in the task of data

Manuscript submitted to ACM

extraction and found that they were in a relatively primitive state in the way they shared raw data and managed files. We believe that online platforms are very necessary in the current situation affected by the COVID-19 pandemic. The online system also allows researchers to work without geographical constraints. DeepShovel uses shared projects and files online to maximize the exposure of team internal data to all team members. Under the premise of this openness, we use the file lock mechanism and the file principal mechanism to ensure the consistency of operations during data extraction.

6.2.4 Cross-team Collaboration Experience. We also studied the case of cross-team collaboration and supported this collaboration in DeepShovel. Our design of user personas gives team leaders control over different members. In the case of cross-team collaboration, the team leader can control the permissions of team members and the security of data to the greatest extent, and users can also establish a new team within the system to support cross-team cooperation. However, we still have users reporting that they want their data to remain private until the end of the study, so we put some restrictions on collaborating with cross-teams. With the emphasis on open science [13], this is still an issue that needs to be considered. When dealing with cross-team collaboration, how to balance the sharing of information with the confidentiality of scientific research can be an important issue.

7 LIMITATIONS AND FUTURE WORK

Our user research and system function design are both conducted in the field of geoscience, especially in the DDE program. We propose a novel and general collaborative framework for scientific literature data extraction in natural science. However, when extended into other disciplines, there may be more problems and difficulties that we still do not fully understand, including the specific data extraction tasks and the form of team collaboration. We plan to expand DeepShovel to other disciplines in the future as more extensive user research is required.

Besides, since our research is conducted remotely through an online meeting service, the setting of such scene may be quite different from the actual situation. We may lack observation of the working status of their team cooperation. For example, we cannot observe how the team communicates in their actual state, how the data is transmitted throughout the teamwork process, and how their data integration process is actually accomplished. Moreover, We do not conduct a comparative experiment study by comparing DeepShovel with some baselines, including extracting data manually and using existing tools like ABBYY. Since the effect of teamwork and the efficiency of data extraction requires long-term observation to obtain, we think we need to observe the team's feedback after long-term use and conduct a larger-scale field deployment.

The rest are some technical limitations in DeepShovel: 1) the quality of data extraction is greatly affected by the quality of PDF files, and we cannot handle some low-resolution scans that are too old; 2) most AI models in DeepShovel are based on rules provided by geoscientists and a relatively small amount of geoscience data, which may lead to some problems in the processing of uncovered literature; 3) in the current proof-of-concept stage, DeepShovel does not meet all the types of data demands (e.g., points location in some scatterplots) because of the lack of relevant datasets. Fortunately, such problems exist in different independent data extraction modules, not affecting the system design framework. In the future, we will continuously add new modules and improve existing modules through rapid system iterative upgrades.

8 CONCLUSION

In this paper, we present DeepShovel, an online collaborative platform for data extraction in the scientific literature with AI assistance that can help researchers cooperate with their teammates to extract data from PDF documents and build a scientific database. The design of DeepShovel is motivated by the user research we place in the field of geoscience. DeepShovel can help researchers extract and aggregate data containing meta information, tables, texts, and location from the literature. The research team can collaborate in DeepShovel, and team members can share resources and progress with others. DeepShovel has been deployed for one month and there are already 253 users from 36 geoscientist teams within the DDE program use it in a daily basis. More than 240 projects and 46,000 documents are being processed for building scientific databases. The follow-up user evaluation with 14 researchers confirms that DeepShovel improves researchers' efficiency in data extraction from geoscience literature and promotes the close collaboration of their teams.

ACKNOWLEDGMENTS

We owe a particular debt of gratitude to the scientists from the Deep-time Digital Earth project who all contributed enormously valuable feedback. We also thank Jia Guo, Yifei Shen, Qi Li, Zhixin Guo, Mingxuan Yan, Mingze Li, Le Zhou, Jingyao Tang, Han Liu, Shengling Zhu and Tao Shi for their support to our system development. This work is supported by National Natural Science Foundation of China (No.42050105, No.62106141) and Shanghai Sailing Program (21YF1421900).

REFERENCES

- [1] 2008–2021. GROBID. <https://github.com/kermitt2/grobid>. arXiv:1:dir:dab86b296e3c3216e2241968f0d63b68e8209d3c
- [2] ABBYY. 1994-2022. PDF software: Open, read & edit pdfs. <https://pdf.abbyy.com/>
- [3] Md Altaf-Ul-Amin, Farit Mochamad Afendi, Samuel Kuria Kiboi, and Shigehiko Kanaya. 2014. Systems biology in the context of big data and networks. *BioMed research international* 2014 (2014).
- [4] Saleema Amershi, Dan Weld, Mihaela Vorvoreanu, Adam Fourney, Besmira Nushi, Penny Collisson, Jina Suh, Shamsi Iqbal, Paul N. Bennett, Kori Inkpen, Jaime Teevan, Ruth Kikin-Gil, and Eric Horvitz. 2019. *Guidelines for Human-AI Interaction*. Association for Computing Machinery, New York, NY, USA, 1–13. <https://doi.org/10.1145/3290605.3300233>
- [5] Zahra Ashktorab, Michael Desmond, Josh Andres, Michael Muller, Narendra Nath Joshi, Michelle Brachman, Aabhas Sharma, Kristina Brimijoin, Qian Pan, Christine T. Wolf, Evelyn Duesterwald, Casey Dugan, Werner Geyer, and Darrell Reimer. 2021. AI-Assisted Human Labeling: Batching for Efficiency without Overreliance. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW1, Article 89 (apr 2021), 27 pages. <https://doi.org/10.1145/3449163>
- [6] Christopher Austin and Fred Kusumoto. 2016. The application of Big Data in medicine: current implications and future directions. *Journal of Interventional Cardiac Electrophysiology* 47, 1 (2016), 51–59.
- [7] Karianne J Bergen, Paul A Johnson, V Maarten, and Gregory C Beroza. 2019. Machine learning for data-driven discovery in solid Earth geoscience. *Science* 363, 6433 (2019).
- [8] Carrie J. Cai, Samantha Winter, David Steiner, Lauren Wilcox, and Michael Terry. 2019. "Hello AI": Uncovering the Onboarding Needs of Medical Practitioners for Human-AI Collaborative Decision-Making. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW, Article 104 (nov 2019), 24 pages. <https://doi.org/10.1145/3359206>
- [9] C Cervato, G Bohling, C Loepp, T Taylor, WS Snyder, P Diver, J Reed, D Fils, D Greer, and Xiaoyun Tang. 2005. The CHRONOS System: geoinformatics for sedimentary geology and paleobiology. In *2005 IEEE International Symposium on Mass Storage Systems and Technology*. IEEE, 182–186.
- [10] Quan Ze Chen, Daniel S Weld, and Amy X Zhang. 2021. Goldilocks: Consistent Crowdsourced Scalar Annotations with Relative Uncertainty. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW2 (2021), 1–25.
- [11] Hao-Fei Cheng, Ruotong Wang, Zheng Zhang, Fiona O'Connell, Terrance Gray, F. Maxwell Harper, and Haiyi Zhu. 2019. *Explaining Decision-Making Algorithms through UI: Strategies to Help Non-Expert Stakeholders*. Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3290605.3300789>
- [12] Christopher Clark and Santosh Divvala. 2016. PDFFigures 2.0: Mining figures from research papers. In *2016 IEEE/ACM Joint Conference on Digital Libraries (JCDL)*. IEEE, 143–152.
- [13] Johanna Cohoon. 2021. *Negotiating Open Science: The Open Science Framework as a Technology-in-Practice*. Association for Computing Machinery, New York, NY, USA, 245–248. <https://doi.org/10.1145/3462204.3481785>

- [14] National Research Council et al. 2000. *A question of balance: Private rights and the public interest in scientific and technical databases*. National Academies Press.
- [15] VI Davyдов and WS Snyder. 2006. Why detail is important: Building a data system for deep-time paleoclimate research. In *AGU Fall Meeting Abstracts*, Vol. 2006. IN53C-05.
- [16] Michael Desmond, Michael Muller, Zahra Ashktorab, Casey Dugan, Evelyn Duesterwald, Kristina Brimijoin, Catherine Finegan-Dollak, Michelle Brachman, Aabhas Sharma, Narendra Nath Joshi, and Qian Pan. 2021. *Increasing the Speed and Accuracy of Data Labeling Through an AI Assisted Interface*. Association for Computing Machinery, New York, NY, USA, 392–401. <https://doi.org/10.1145/3397481.3450698>
- [17] Rodolfo Dirzo, Hillary S Young, Mauro Galetti, Gerardo Ceballos, Nick JB Isaac, and Ben Collen. 2014. Defaunation in the Anthropocene. *science* 345, 6195 (2014), 401–406.
- [18] Jun-xuan Fan, Shu-zhong Shen, Douglas H Erwin, Peter M Sadler, Norman MacLeod, Qiu-ming Cheng, Xu-dong Hou, Jiao Yang, Xiang-dong Wang, Yue Wang, et al. 2020. A high-resolution summary of Cambrian to Early Triassic marine invertebrate biodiversity. *Science* 367, 6475 (2020), 272–277.
- [19] Thomas A Finholt and Gary M Olson. 1997. From laboratories to collaboratories: A new organizational form for scientific collaboration. *Psychological Science* 8, 1 (1997), 28–36.
- [20] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep learning*. MIT press.
- [21] Vidhya Govindaraju, Ce Zhang, and Christopher Ré. 2013. Understanding tables in context using standard NLP toolkits. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 658–664.
- [22] Gótz Hoeppe. 2021. Encoding Collective Knowledge, Instructing Data Reusers: The Collaborative Fixation of a Digital Scientific Data Set. *Computer Supported Cooperative Work (CSCW)* 30, 4 (2021), 463–505.
- [23] Zhi Hong, Logan Ward, Kyle Chard, Ben Blaiszik, and Ian Foster. 2021. Challenges and Advances in Information Extraction from Scientific Literature: a Review. *JOM* 73, 11 (2021), 3383–3400.
- [24] Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. (2017). To appear.
- [25] Eric Horvitz. 1999. Principles of Mixed-Initiative User Interfaces. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '99)*. Association for Computing Machinery, New York, NY, USA, 159–166. <https://doi.org/10.1145/302979.303030>
- [26] Anuj Karpatne, Imme Ebert-Uphoff, Sai Ravela, Hassan Ali Babaie, and Vipin Kumar. 2018. Machine learning for the geosciences: Challenges and opportunities. *IEEE Transactions on Knowledge and Data Engineering* 31, 8 (2018), 1544–1554.
- [27] Anthony Kay. 2007. Tesseract: An Open-Source Optical Character Recognition Engine. *Linux J.* 2007, 159 (jul 2007), 2.
- [28] Candice Lanius, Kerstin Lehner, and Lucia Profeta. 2021. Usability Testing EarthChem PetDB: Findings for Interface Improvement and Recommendations for Future User Testing. In *AGU Fall Meeting 2021*. AGU.
- [29] Chen Li, Maria Liakata, and Dietrich Rebholz-Schuhmann. 2014. Biological network extraction from scientific literature: state of the art and challenges. *Briefings in bioinformatics* 15, 5 (2014), 856–877.
- [30] Minghao Li, Lei Cui, Shaohan Huang, Furu Wei, Ming Zhou, and Zhoujun Li. 2019. TableBank: A Benchmark Dataset for Table Detection and Recognition. [arXiv:cs/1903.01949](https://arxiv.org/abs/cs/1903.01949)
- [31] Robyn Longhurst. 2003. Semi-structured interviews and focus groups. *Key methods in geography* 3, 2 (2003), 143–156.
- [32] Yi Luan. 2018. Information extraction from scientific literature for method recommendation. *arXiv preprint arXiv:1901.00401* (2018).
- [33] William J McMahon and Neil S Davies. 2018. Evolution of alluvial mudrock forced by early land plants. *Science* 359, 6379 (2018), 1022–1024.
- [34] Hendrik Müller, Aaron Sedley, and Elizabeth Ferrall-Nunge. 2014. Survey research in HCI. In *Ways of Knowing in HCI*. Springer, 229–266.
- [35] Feng Niu, Che Zhang, Christopher Ré, and Jude W Shavlik. 2012. DeepDive: Web-scale Knowledge-base Construction using Statistical Learning and Inference. *VLDS* 12 (2012), 25–28.
- [36] Roland Oberhänsli. 2020. Deep-time Digital Earth (DDE) the First IUGS Big Science Program. *Journal of the Geological Society of India* 95, 3 (2020), 223–226.
- [37] Judith S Olson and Wendy A Kellogg. 2014. *Ways of Knowing in HCI*. Vol. 2. Springer.
- [38] Shanan E Peters, Jon M Husson, and Julia Wilcots. 2017. The rise and fall of stromatolites in shallow marine environments. *Geology* 45, 6 (2017), 487–490.
- [39] Stephen J. Puetz. 2018. A relational database of global U-Pb ages. *Geoscience Frontiers* 9, 3 (2018), 877–891. <https://doi.org/10.1016/j.gsf.2017.12.004>
- [40] Stephen J Puetz, Carlos E Ganade, Udo Zimmermann, and Glenn Borchardt. 2018. Statistical analyses of global U-Pb database 2017. *Geoscience Frontiers* 9, 1 (2018), 121–145.
- [41] Kohulan Rajan, Henning Otto Brinkhaus, Maria Sorokina, Achim Zielesny, and Christoph Steinbeck. 2021. DECIMER-Segmentation: Automated extraction of chemical structure depictions from scientific literature. *Journal of Cheminformatics* 13, 1 (2021), 1–9.
- [42] Alexander Ratner, Stephen H. Bach, Henry R. Ehrenberg, Jason Alan Fries, Sen Wu, and Christopher Ré. 2017. Snorkel: Rapid Training Data Creation with Weak Supervision. *CoRR* abs/1711.10160 (2017). arXiv:1711.10160 <http://arxiv.org/abs/1711.10160>
- [43] Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh, and Lora M Aroyo. 2021. “Everyone Wants to Do the Model Work, Not the Data Work”: Data Cascades in High-Stakes AI. Association for Computing Machinery, New York, NY, USA. <https://doi.org/10.1145/3411764.3445518>

- [44] Gabriel Ravanhani Schleder, Antonio Claudio M Padilha, Alexandre Reily Rocha, Gustavo Martini Dalpian, and Adalberto Fazzio. 2019. Ab initio simulations and materials chemistry in the age of big data. *Journal of chemical information and modeling* 60, 2 (2019), 452–459.
- [45] Kjeld Schmidt. 2008. Taking CSCW Seriously: Supporting Articulation Work (1992). In *Cooperative Work and Coordinative Practices*. Springer, 45–71.
- [46] Peters Shanan, Livny Miron, Rekatsinas Theo, Ross Ian, Quinn Daven, Glassel Aimee, Aydemir Brian, Edquist Carl, Peterson Jeff, Ré Christopher, Husson Jon, Zaffos Andrew, Wilcots Julia, Czaplewski John, Syverson Valerie, Zhang Ce, Ito Erika, Liu Chao, Wieferich Daniel, and Serna Brandon. [n. d.]. xDD: A digital library and cyberinfrastructure facilitating the discovery and utilization of data & knowledge in published documents. <https://xdd.wisc.edu/about.html>.
- [47] WS Snyder, KA Lehnert, E Ito, Ulrich Harms, and Jens Klump. 2008. GeosciNET: Building a Global Geoinformatics Partnership. In *AGU Fall Meeting Abstracts*, Vol. 2008. IN31D-03.
- [48] Stephanie B Steinhardt and Steven J Jackson. 2014. Reconciling rhythms: plans and temporal alignment in collaborative scientific work. In *Proceedings of the 17th ACM conference on Computer supported cooperative work & social computing*. 134–145.
- [49] Ziheng Sun, Laura Sandoval, Robert Crystal-Ornelas, S Mostafa Mousavi, Jinbo Wang, Cindy Lin, Nicoleta Cristea, Daniel Tong, Wendy Hawley Carande, Xiaogang Ma, et al. 2022. A review of Earth Artificial Intelligence. *Computers & Geosciences* (2022), 105034.
- [50] Matthew C Swain and Jacqueline M Cole. 2016. ChemDataExtractor: a toolkit for automated extraction of chemical information from the scientific literature. *Journal of chemical information and modeling* 56, 10 (2016), 1894–1904.
- [51] Dominika Tkaczyk, Andrew Collins, Paraic Sheridan, and Joeran Beel. 2018. Machine Learning vs. Rules and Out-of-the-Box vs. Retrained: An Evaluation of Open-Source Bibliographic Reference and Citation Parsers. arXiv:cs.DL/1802.01168
- [52] Dominika Tkaczyk, Paweł Szostek, Piotr Jan Dendek, Mateusz Fedoryszak, and Łukasz Bolikowski. 2014. Cermine—automatic extraction of metadata and references from scientific literature. In *2014 11th IAPR International Workshop on Document Analysis Systems*. IEEE, 217–221.
- [53] Dominika Tkaczyk, Paweł Szostek, Mateusz Fedoryszak, Piotr Jan Dendek, and Łukasz Bolikowski. 2015. CERMINE: automatic extraction of structured metadata from scientific literature. *International Journal on Document Analysis and Recognition (IJDAR)* 18, 4 (2015), 317–335.
- [54] Marlee A Tucker, Katrin Böhning-Gaese, William F Fagan, John M Fryxell, Bram Van Moorter, Susan C Alberts, Abdullahi H Ali, Andrew M Allen, Nina Attias, Tal Avgar, et al. 2018. Moving in the Anthropocene: Global reductions in terrestrial mammalian movements. *Science* 359, 6374 (2018), 466–469.
- [55] Theresa Velden. 2013. Explaining field differences in openness and sharing in scientific communities. In *Proceedings of the 2013 conference on Computer supported cooperative work*. 445–458.
- [56] Vickie R Walker, Charles P Schmitt, Mary S Wolfe, Artur J Nowak, Kuba Kulesza, Ashley R Williams, Rob Shin, Jonathan Cohen, Dave Burch, Matthew D Stout, et al. 2022. Evaluation of a semi-automated data extraction tool for public health literature-based reviews: Dextr. *Environment international* 159 (2022), 107025.
- [57] Chengshan Wang, Robert M Hazen, Qiuming Cheng, Michael H Stephenson, Chenghu Zhou, Peter Fox, Shu-zhong Shen, Roland Oberhänsli, Zengqian Hou, Xiaogang Ma, Zhiqiang Feng, Junxuan Fan, Chao Ma, Xiumian Hu, Bin Luo, Juanle Wang, and Craig M Schiffries. 2021. The Deep-Time Digital Earth program: data-driven discovery in geosciences. *National Science Review* 8, 9 (02 2021). <https://doi.org/10.1093/nsr/nwab027>
- [58] Nancy Xin Ru Wang, Douglas Burdick, and Yunyao Li. 2021. TableLab: An Interactive Table Extraction System with Adaptive Deep Learning. In *26th International Conference on Intelligent User Interfaces*. 87–89.
- [59] Mark D Wilkinson, Michel Dumontier, Ijsbrand Jan Aalbersberg, Gabrielle Appleton, Myles Axton, Arie Baak, Niklas Blomberg, Jan-Willem Botten, Luiz Bonino da Silva Santos, Philip E Bourne, et al. 2016. The FAIR Guiding Principles for scientific data management and stewardship. *Scientific data* 3, 1 (2016), 1–9.
- [60] Sen Wu, Luke Hsiao, Xiao Cheng, Braden Hancock, Theodoros Rekatsinas, Philip Levis, and Christopher Ré. 2018. Fonduer: Knowledge Base Construction from Richly Formatted Data. In *Proceedings of the 2018 International Conference on Management of Data (SIGMOD '18)*. Association for Computing Machinery, New York, NY, USA, 1301–1316. <https://doi.org/10.1145/3183713.3183729>
- [61] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. 2019. Detectron2. <https://github.com/facebookresearch/detectron2>.
- [62] Ce Zhang. 2015. *DeepDive: a data management system for automatic knowledge base construction*. Ph.D. Dissertation. The University of Wisconsin-Madison.
- [63] Ce Zhang, Vidhya Govindaraju, Jackson Borchardt, Tim Foltz, Christopher Ré, and Shanan Peters. 2013. GeoDeepDive: Statistical Inference Using Familiar Data-Processing Languages. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data (SIGMOD '13)*. Association for Computing Machinery, New York, NY, USA, 993–996. <https://doi.org/10.1145/2463676.2463680>