# Accelerating Geoscience Research: An Advanced Platform for Efficient Multimodal Data Integration from Geoscience Literature

**Zhixin Guo[1], Jianping Zhou[1], Guanjie Zheng[2], Xinbing Wang[1], and Chenghu Zhou[3]**

[1]Department of Electronic Engineering, Shanghai Jiao Tong University (SJTU), [2]John Hopcroft Center for Computer Science, SJTU, [3]Institute of Geographic Sciences and Natural Resources Research, CAS

## Abstract

In the era of big data science, geoscience has experienced a significant paradigm shift, moving towards a data-driven approach to scientific discovery. This shift, however, presents a considerable challenge due to the dispersed nature of data. To address this issue, we introduce a comprehensive, publicly accessible platform designed to facilitate the extraction and integration of multimodal data from geoscience literature, encompassing text, visual, and tabular formats. Our platform has been effectively applied in processing diverse domains, including mountain disaster data, global orogenic belt isotope data, and environmental pollutant data. This has facilitated substantial academic research, evidenced by developing **knowledge graphs based on debris flow data**, establishing a **global Sm-Nd isotope database**, and meticulous detection and analysis of environmental pollutants. As a core component of the Deep-Time Digital Earth (DDE) program, our platform has significantly contributed to the field, supporting forty geoscience research teams in their endeavors and processing over 40,000 documents. This accomplishment underscores the platform's capacity for handling large-scale data and its pivotal role in advancing geoscience research in the age of big data.

(@ Jianping Zhou)   (@ GeoKnowledgeFusion)

## Background & Motivation

- Geological data exists in a multimodal form in a large number of academic literature, including tables, images, text and other forms.
- Geoscientific research in the era of big data often needs to mine knowledge from large amounts of geoscientific data and synthesize the correlations between them.
- However, it is very difficult to extract and normalize the heterogeneous data from multiple sources.
- Therefore, we hope to design a one-stop standardized process platform from literature to geoscientific knowledge to accelerate the research of geoscientists.

## Our method & System

- a comprehensive, publicly accessible platform designed to facilitate the extraction and integration of multimodal data from geoscience literature.
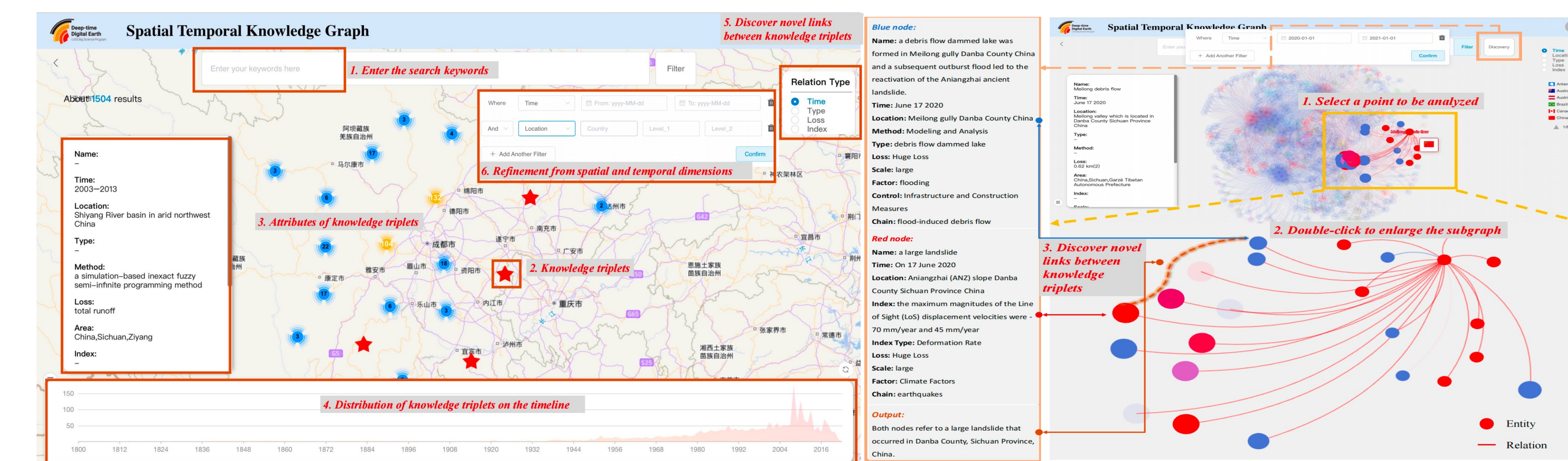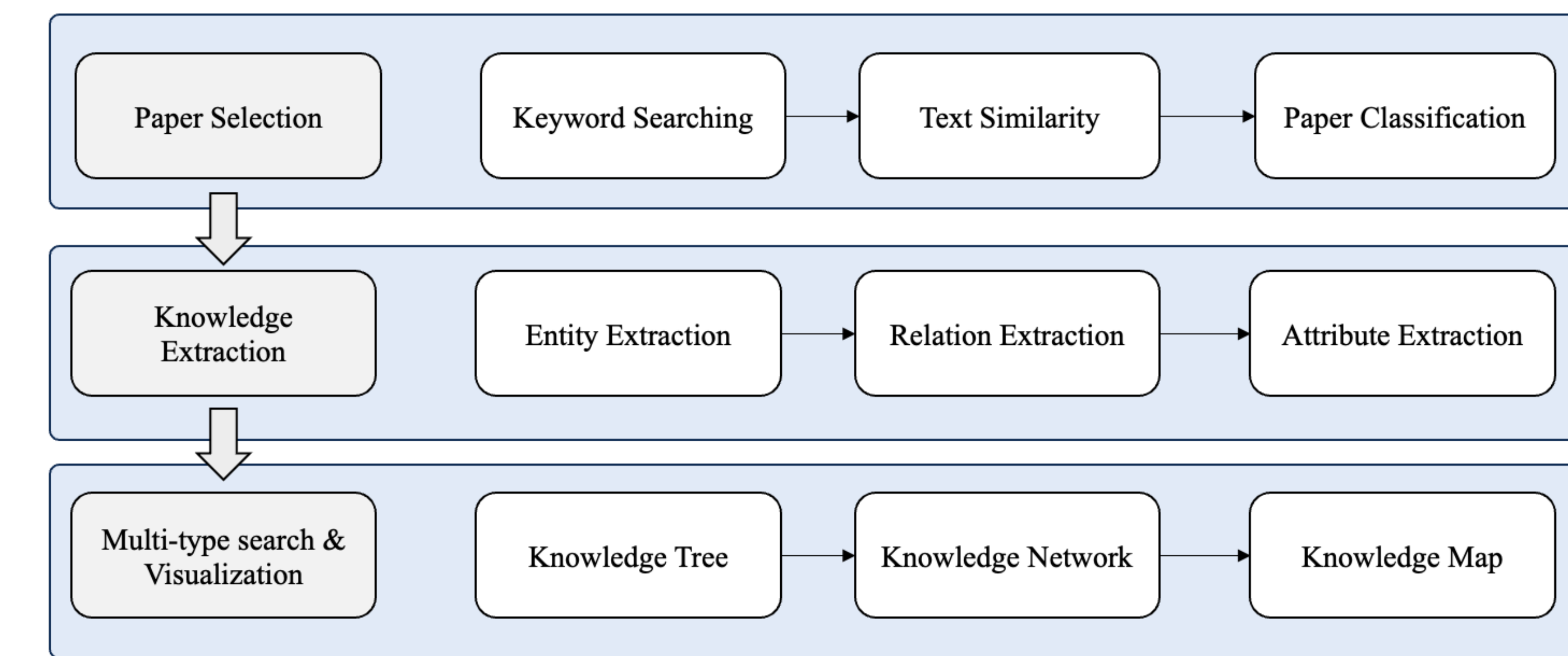


## Case 1: The Global Sm-Nd Isotope Tabular Database

- Sm-Nd isotope data are widely available in geoscientific literature tables, necessitating the efficient integration of this data into a database through automated extraction and fusion method.
- We collected 270,527 academic papers from 33 types of geology journals, screened 10,548 Sm-Nd isotope data, parsed 42,432 tables, screened 4,807 valid tables, and finally parsed 17,334 valid data.



## Case 2: Debris Flow Knowledge Graph

- The global debris flow lacks systematic and refined information on the details of the hazard, and the mechanism that causes the hazard, which are widely available in the texts of geoscience literature.
- We screened 55,360 articles of debris flow literature from 136 types of mountain disaster journals, and extracted 18 types of debris flow hazard indicators.



## Related Paper

[1] Zhixin Guo, Jianping Zhou, et al., "Towards Controlled Table-to-Text Generation with Scientific Reasoning", ICASSP, 2024.

[2] Zhixin Guo, et al. "Sm-Nd Isotope Data Compilation from Geoscientific Literature Using an Automated Tabular Extraction Method", 2024.