

# OneDZ: A Global Detrital Zircon Database and Implications for Constructing Giant Geoscience Database

Keran Li<sup>1</sup>, Xiumian Hu<sup>1</sup>\*, Rong Chai<sup>2</sup>, Jianghai Yang<sup>3</sup>, Weiwei Xue<sup>4</sup>, Yingdi Pan<sup>1</sup>, Taiyang Li<sup>5</sup>, Can Fang<sup>6</sup>, Anlin Ma<sup>1</sup>, Hu Huang<sup>7,8</sup>, Qianqian Guo<sup>9</sup>, Wentao Yang<sup>10</sup>, Lisha Hu<sup>11</sup>, Liang Qi<sup>7,8</sup>, Guohui Chen<sup>12</sup>,  
5 Gaoyuan Sun<sup>13</sup>, Shijie Zhang<sup>14</sup>, Tao Deng<sup>1</sup>, Kuizhou Li<sup>7,15</sup>, Jiaopeng Sun<sup>16</sup>, Biao Gao<sup>17</sup>

<sup>1</sup>State Key Laboratory of Mineral Deposit Research, School of Earth Sciences and Engineering, Nanjing University, Nanjing 210023, China

<sup>2</sup>Chinese Academy of Geological Sciences, Beijing 100037, China

<sup>3</sup>School of Earth Sciences, China University of Geosciences (Wuhan), Wuhan 430074, China

10 <sup>4</sup>State Key Laboratory of Isotope Geochemistry, Guangzhou Institute of Geochemistry, Chinese Academy of Sciences (CAS), Guangzhou, China

<sup>5</sup>College of Computer Science and Cyber Security, Chengdu University of Technology, Chengdu 610059, China

<sup>6</sup>Hangzhou Research Institute, Huawei Technologies, Hangzhou 310056, China

15 <sup>7</sup>State Key Laboratory of Oil and Gas Reservoir Geology and Exploitation, Institute of Sedimentary Geology, Chengdu University of Technology, Chengdu 610059, China

<sup>8</sup>Key Laboratory of Deep-time Geography and Environment Reconstruction and Applications of Ministry of Natural Resources, Chengdu University of Technology, Chengdu 610059, China

<sup>9</sup>College of Earth and Planetary Sciences, University of Chinese Academy of Sciences, Beijing, 100049, China

<sup>10</sup>School of Resources and Environment, Henan Polytechnic University, Jiaozuo 454000, China

20 <sup>11</sup>College of Marine Geosciences, Ocean University of China, Qingdao 266100, China

<sup>12</sup>School of Earth Sciences and Engineering, Hohai University, Nanjing 210098, China

<sup>13</sup>College of Oceanography, Hohai University, Nanjing 210024, China

<sup>14</sup>College of Tourism, Henan Normal University, Xinxiang, China

<sup>15</sup>College of Earth and Planetary Sciences, Chengdu University of Technology, Chengdu 610059, China

25 <sup>16</sup>State Key Laboratory of Continental Dynamics, Department of Geology, Northwest University, Xi'an 710069, China

<sup>17</sup>State Key Laboratory of Palaeobiology and Stratigraphy, Nanjing Institute of Geology and Palaeontology and Center for Excellence in Life and Palaeoenvironment, Chinese Academy of Sciences, Nanjing 210008, China

\* *Correspondence to:* Xiumian Hu (huxm@nju.edu.cn)

**Abstract.** The amount of detrital zircon U-Pb geochronology data and Lu-Hf isotopic data has doubled in the past two decades  
30 with the continuous improvement of analytical methods, and has developed into the most closely integrated research field in  
earth science with big data methods. However, how to effectively construct giant databases in geoscience has become a  
challenge. Here, we present OneDZ, a global comprehensive detrital zircon U-Pb geochronology and Lu-Hf isotope database,  
which includes diverse samples with data source, location, stratigraphy, depositional age, and various elemental and isotopic  
information. OneDZ collected corresponding regions, stratigraphic and lithological information to facilitate quick access for  
35 users. Comparing with current zircon databases, Till now, OneDZ complies 2,550,738 grains of detrital zircon U-Pb and  
297,527 grains of detrital zircon Lu-Hf records from 275,971 publications. Furthermore, the construction of OneDZ leverages  
artificial intelligence (AI) and programming scripts and offers insights into managing large-scale unstructured data in  
geosciences. This paper further discusses the perspective of applying big data methods in the research of zircon-related areas.

This database exemplifies the power of big data in Earth sciences, providing a platform for investigating zircon data in deep  
40 time. It serves as a springboard for research, offering new insights in understanding Earth's past, present, and future. The  
database (Li et al., 2026) is freely available via Zenodo at 10.5281/zenodo.19690702.

## 1 Introduction

The advent of high-precision U-Pb geochronology has revolutionized our understanding of Earth history. Isotope-dilution  
thermal-ionization mass spectrometry (ID-TIMS; Krogh, 1973) remains the benchmark for highest accuracy and precision  
45 ( $\leq 0.1\%$ ), but its destructive, time-intensive protocol limits statistical throughput for large detrital suites. Laser-ablation  
inductively-coupled-plasma mass spectrometry (LA-ICP-MS) and secondary-ion mass spectrometry (SIMS) now provide  
rapid, in-situ analyses with 1 - 3 % precision, which is ideal for analyzing the hundreds to thousands of concordant ages  
required for robust detrital-zircon provenance and maximum depositional-age studies (e.g., Jarvis and Kym, 1988; MacRae  
and Neil, 1995; Belu et al., 2003; Muzikar et al., 2003; Yergey et al., 2013). Together, these complementary techniques extend  
50 high-precision geochronology from single crystals to entire sedimentary systems. Zircon, a robust and ubiquitous mineral  
found throughout the continental crust, serves as a reliable recorder of geological events due to its high closure temperature  
and resistance to weathering and metamorphism (Pupin, 1980). Primary zircons from magmatic or metamorphic rocks are  
commonly fragmented, transported, and ultimately deposited as detrital zircons.

Typical analyses of detrital zircons include U-Pb and Lu-Hf isotopic systems. Chemical formula of detrital zircon can be  
55 represented as  $[\text{ZrSiO}_4]$ . The ionic radius of  $[\text{Zr}^{4+}]$  is  $0.87 \text{ \AA}$ , which can be easily replaced by  $[\text{U}^{4+}]$  and  $[\text{Th}^{4+}]$  because of  
similar ionic radius of  $1.05 \text{ \AA}$  and  $1.10 \text{ \AA}$  (Jaffey et al., 1971). Two isotopes of  $[\text{U}^{4+}]$  ( $^{238}\text{U}$  and  $^{235}\text{U}$ ) generate  $^{206}\text{Pb}$  and  $^{207}\text{Pb}$   
isotopes following the decay processes:  $^{238}\text{U} \rightarrow ^{206}\text{Pb} + 8\alpha + 6\beta^-$  (half-life: 4468 million years, Jaffey et al., 1971),  
 $^{235}\text{U} \rightarrow ^{207}\text{Pb} + 7\alpha + 4\beta^-$  (half-life: 703.8 million years, Jaffey et al., 1971) and  $^{232}\text{Th} \rightarrow ^{208}\text{Pb} + 6\alpha + 4\beta^-$  (half-life: 1400 million years,  
Jaffey et al., 1971). Based on triple decay processes, the detrital zircon ages can be obtained via consisted  $^{206}\text{Pb}/^{238}\text{U}$ ,  $^{207}\text{Pb}/^{235}\text{U}$   
60 and  $^{208}\text{Pb}/^{232}\text{Th}$  decayed ages.

In addition to U-Pb geochronology, the Lu-Hf isotopic system has become an indispensable tool for understanding crustal  
evolution and mantle differentiation (Patchett and Tatsumoto, 1983). The  $[\text{Lu}^{3+}]$  is the heaviest rare earth element (REE) and  
are easily enriched in detrital zircon. The  $^{176}\text{Lu}$  decays to  $^{176}\text{Hf}$  via  $^{176}\text{Lu} \rightarrow ^{176}\text{Hf} + \beta^-$  (half-life: 37.1 billion years, Kinny and  
Mass, 2003). Except for the geochronological application, the Lu-Hf isotopic data can be used to gain the original information  
65 (Scherer et al., 2001; Söderlund et al., 2004). The Lu-Hf isotopic data are noted by  $\epsilon$  units by  
 $\epsilon_{\text{Hf}}(0) = 10000 \times \left[ \frac{(^{176}\text{Hf}/^{177}\text{Hf})_{\text{sample}}}{(^{176}\text{Hf}/^{177}\text{Hf})_{\text{CHUR},0}} - 1 \right]$  and  $\epsilon_{\text{Hf}}(t) = 10000 \times \left\{ \left[ \frac{(^{176}\text{Hf}/^{177}\text{Hf})_{\text{sample}} - (^{176}\text{Lu}/^{177}\text{Hf})_{\text{sample}} \times (e^{\lambda t} - 1)}{(^{176}\text{Hf}/^{177}\text{Hf})_{\text{CHUR},0} - (^{176}\text{Lu}/^{177}\text{Hf})_{\text{CHUR}} \times (e^{\lambda t} - 1)} \right] - 1 \right\}$ .  $t$  is the crystallization age.  $^{176}\text{Hf}/^{177}\text{Hf}$  and  $^{176}\text{Lu}/^{177}\text{Hf}$  can be measured  
from detrital zircons.  $\lambda$  is the decay constant and equals to  $1.867 \times 10^{-5}$  million years (Söderlund et al., 2004). CHUR denotes  
the isotopic results of the chondritic uniform reservoir.

70 In the past two decades, it is estimated that millions U-Pb geochronological data of detrital zircons have been internationally reported. As the amount of data increases, it's possible to use detrital zircon data with big data methods for analyzing significant scientific problems. For instance, the compilation of detrital zircon big data is used for the reconstruction of continental arcs (McKenzie et al., 2016; Cao et al., 2017), tectonic history (Cawood et al., 2012; Barham et al., 2022; Zhang et al., 2023; Malone et al., 2024; Odlum et al., 2024), crustal evolution (Cheng, 2017; Barham et al., 2019; Cawood, 2020), paleo-  
75 geographic (Xue et al., 2022; Jian et al., 2022) and provenance analysis (Wang et al., 2024). Along with data-driven analysis, several analytical tools (Ludwing, 2003; Vermeesch, 2018; Saylor et al., 2017; Sharman et al., 2018) and professional databases have been established (Voice et al., 2011; Puetz, 2019; Martin et al., 2022; Puetz, 2024; Wu et al., 2024). However, the existing databases are not primarily designed for the needs of sedimentological researches and the reported data are usually mixed with magmatic and metamorphic rocks. With the rapid accumulation of detrital zircon data, existing databases are difficult to  
80 effectively cover detrital zircon data in sedimentary rocks.

Current database compilations in the Earth sciences predominantly emphasize data dissemination while offering limited discussion of the procedural challenges inherent to data construction. Previous databases typically report compiled datasets directly while neglecting the complexities of data collection and processing (Voice et al., 2011; Puetz, 2019; Martin et al., 2022; Puetz, 2024; Wu et al., 2024; Table 1), an omission that has caused database growth to lag substantially behind the rate  
85 of scientific publication. Consequently, the construction of large-scale geoscience databases, particularly for detrital zircon data, urgently necessitates a shift from manual curation toward systematic, automated collection methodologies, especially as established repositories such as EarthChem, GEOROC, EarthBank, and Geochron, despite providing user-friendly web interfaces for data querying and submission, still require contributors to manually extract and reformat data from original publications into standardized templates, a time-consuming bottleneck that continues to impede data contribution. To address  
90 these challenges, we established OneDZ, a comprehensive database of detrital zircon U-Pb geochronological and Lu-Hf isotopic data covering global English and Chinese literature through 2022. Inspired by the emerging "literature-as-datasets" paradigm utilizing large language models (LLMs), we experimentally deployed multiple automated LLM-driven agents for data collection, enabling users to contribute through original PDF files or DOI information alone.

Currently OneDZ encompasses 2,550,738 U-Pb and 297,527 Lu-Hf records from approximately 275,971 publications.  
95 Furthermore, we implemented a dual-track quality assurance system in which automated agent-based extraction and verification facilitate rapid data proliferation, while expert inspection ensures reliability, with 1,414,062 U-Pb records and all Lu-Hf records having been manually verified by specialists. Data in OneDZ spans nearly the entire history of earth's sediments, offering valuable insights into the timing and nature of geological events. The compilation includes data from various analytical techniques, host rock lithologies, stratigraphic information, and other original records. OneDZ records the lithology,  
100 stratigraphic, spatial, and testing information of detrital zircons as much as possible. In the compilation of OneDZ, Python scripts were developed for systematic data cleaning, format standardisation, and quality control, ensuring that the final dataset is internally consistent and ready for immediate reuse. The enormous volume of data also makes OneDZ a valuable resource

for discussing data analysis methods in Earth science. OneDZ provides a foundation for research in multiple aspects, including data provision, data harmonisation, and discussion and analysis of data analysis methods in earth science.

## 105 2 Database construction

One of the most unique features in OneDZ is the systematic construction workflow (Figure 1). Firstly, the knowledge graph (Hu et al., 2024) was adopted and guided the header design by identifying the most frequent words related to detrital zircon. With the knowledge graph, the words to describe the sample location, sedimentary or stratigraphic descriptions and the isotopic results are the most relevant information associated with detrital zircon studies. Previous research rarely summarized the difficulties in collecting data sources. Guo et al. (2024) summarized current geoscience data compilation challenges include non-repeatability, uncertainty, multi-dimensionality, computational complexity, and frequent updates, which pose significant obstacles to the efficient collection and management of geological information. In practice, constantly switching potential literature search engines and manually downloading potential articles one by one actually occupies the main time of database construction. In this project, AI-assisted tools, including DataExpo and GPT Agent (Supplement 1) and large language model, were employed to check specific online resources like Pangea (<https://pangaea.de/>), Google Scholar, and CNKI to search potential papers containing data and capture meta data from PDF files. Following the AI tools, manual verification was conducted, and publication information were passed to several volunteering experts based on their interest regions. These experts extracted and cleaned data using the computer-vision tool, DeepShovel (Zhang et al., 2023), and Python/SQL scripts. These validated data were imported into the OneDZ database. Table 1 provides a detailed comparison with other detrital-zircon databases.

### 2.1 Crowdfunded construction

In the era of data explosion, crowdfunding has become an efficient method for building mega databases. Inspired by this cooperative construction, the OneDZ database was established by dividing different regions and quickly organizing a group of experts in detrital zircons. The crowdfunding approach ensures that each scientist is familiar with the contributed data, maximizing efficiency and accuracy within the same framework following a standard. This method also facilitates dynamic database updates and promotes sustained growth in data volume. The crowdfunded construction is anchored by several regional detrital zircon databases mainly in China which published in a special issue of the journal of Geosciences Data Journal (see Yang et al., 2023), including those from the North China Block (Yang et al., 2023; Dong et al., 2023), the Eastern Central China Orogenic Belt (Chai, 2023), the Songpan-Ganzi and Western Qinling terranes (Pan et al., 2023), the Central Asian Orogenic Belt (Wang, 2023), South China (Luo et al., 2023; Xia, 2023), the Qilian-Qaidam-Kunlun collage (He, 2023), the South China Sea (Huang et al., 2023), the Tarim-West Kunlun-Pamir-Tajik-Tianshuihai terranes (Zhang et al., 2023), the Middle East (Chen et al., 2023; Sun et al., 2023), and samples from Quaternary sediments (Chen et al., 2023). However, the experts-driven crowdfunding could not ensure collecting all data. Therefore, for publications after 2022, a new approach was

135 applied, where everyone can just offer PDF files (or even just the DOI information) and the specific number would be extracted by AI agent. This method lowers the technical threshold for data collecting.

## 2.2 Facility from AI tools

One of the fundamental challenges in compiling large geoscience datasets lies in data collection. Although the crowd-funded approach ensures that geologists participating in database development are experts in their research area, their expertise does not guarantee familiarity with every publication. To find potential metadata, this study introduced a data parsing system integrated with deep learning technology. The data parsing tool is named DataExpo (Lu et al., 2023) and employs deep learning for metadata extraction (Figure S1-S2 in the Supplement), performing automatic semantic tagging, classification, and structured information extraction from web pages. DataExpo automatically crawls web pages related to detrital zircon research. Using a multidimensional web page ranking strategy, retrieval results for different queries are sorted. Finally, based on natural language processing (NLP) and convolutional neural networks (CNNs), DataExpo adjusts the ranking of retrieval results and determines whether to push them to experts. Another AI tool, AI Agent, was created through prompt engineering to analyze characters from specific websites and find potential titles about detrital zircons. Details on using DataExpo and GPT Agent in the OneDZ database construction are provided in the Supplement 1.

In addition to integrating data sources, data extraction poses another major challenge. While most online articles store data in Excel tables as attachments, a considerable number of detrital zircon data is stored in the main text in the article either in table or in text form. To accelerate construction, the interactive computer-vision AI tool DeepShovel (Zhang et al., 2023) was utilized to automatically split tables via optical character recognition. Details on using DeepShovel can be found in Figure S3 in the Supplement. To attract more data contribution and lower the technical requirements, automatic data collecting tool from large-language model agent was applied. After receiving the contributed PDF file, the multi-modal agent automatically extracts the metadata and saves as json files. Then another data-checking agent would automatically check the data quality. At last, an independent agent would evaluate each item. Only items with over 60 % information were provided would be sent to experiencing manual checking. Details on using DeepShovel can be found in Figure S3 in the Supplement.

## 2.3 Automatic data process

In the era of exponential data growth, the construction of domain-specific earth-science databases is becoming the norm. Yet existing zircon and broader geoscience repositories overwhelmingly emphasize data quality, while the critical step of data cleaning has received little systematic attention. In the construction of large scientific databases, beyond ensuring the quality of the original data, it is also essential to trace and maintain the quality of different versions of data formed during the database construction process, a procedure known as data cleaning. Hellerstein et al. (2013) and Chu et al. (2016) identified the key steps in the data cleaning process including (1) Data review and understanding; (2) Missing value processing; (3) Outlier detection and handling; (4) Data format and type conversion; (5) Data consistency and normalization; (6) Data de-duplication. Following the standard data cleaning process, Python scripts were designed for detecting missing key items, checking for

conflicting content, detecting format anomalies, and eliminating duplicate data entries (see Supplement 2 and Figure S4 in the Supplement for details).

### 3 Data compilation and harmonisation

170 The OneDZ dataset is distributed as flat CSV files with a uniform, standardised column schema to maximise interoperability and ease of reuse. Each row represents a single detrital zircon analysis, and columns are organised into thematic groups: bibliographic metadata, sample location and stratigraphy, depositional age constraints, analytical method details, isotopic ratios and their uncertainties, calculated ages, and elemental concentrations (Figure 2). All U-Pb records follow a single 64-column schema, while Lu-Hf records follow a 33-column schema (see README file in the Zenodo repository for the complete field dictionary). This flat-file structure ensures that users can load, filter, and analyse the data with any standard statistical software (e.g., R, Python, Excel, Matlab) without requiring database connectivity or knowledge of relational table structures.

175 In the OneDZ dataset we have compiled 2,550,738 detrital zircon grains (1,414,062 after expert verification) with U-Pb ages and 297,527 grains with Lu-Hf isotope data. From multiple dimensions such as region, literature, and samples, OneDZ is currently the most comprehensive compilation for global detrital zircon data records (Table 2). The U-Pb geochronological data are spatially distributed across 142 geographic regions (Figure 3a). The Lu-Hf data are primarily distributed across China, South Africa, India, and Australia (Figure 3b). Periodic statistics indicate that ancient zircons (over 1000 Ma) predominantly contribute to this dataset in both U-Pb and Lu-Hf data (Figure 3c-f). The content and completeness of sample metadata, spatial data, and stratigraphic information in OneDZ are summarized in Tables 3-6 (expressed as the ratio of valid entries to total entries).

180

#### 3.1 Reference information

185 The reference information in OneDZ includes the principal investigator, publication year, journal name, volume, pagination, article title, and a direct weblink to the original publications. Figure S5 in the Supplement provides a temporal overview of the geographic distribution of these scholarly works. OneDZ aggregates a comprehensive total of 742,832 papers from 1995 to 2022 (Figure S5a Supplement), which includes 52,604 English-language papers and 203,326 Chinese-language papers in the U-Pb datasets. For the Lu-Hf datasets, the compilation consists of 65,420 English-language papers and 8,762 Chinese-language papers from 2004 to 2022 (Figure S5b in the Supplement). Additionally, publicly available master's and doctoral dissertations have been incorporated into the dataset. To ensure accessibility and inclusivity, Chinese-language papers on detrital zircons have been meticulously translated into English. In the U-Pb age dataset, journals such as Precambrian Research, Geological Society of American Bulletin, and Gondwana Research predominantly contribute to the database (Figure S6a-b in the Supplement). The Lu-Hf analyses in OneDZ are drawn entirely from the same journal pool that provided the U-Pb data (Figure S6c-d in the Supplement). Comparing with previous databases (Puetz et al., 2024; Wu et al., 2023), OneDZ surpasses existing repositories in volume and in journal diversity (Figure S6a-b in the Supplement).

195

### 3.2 Sample, spatial and strata information

OneDZ contains the published sample ID, country or state, region, continent, major and minor geographic or geological description of the sediments. In geological research, geological bodies, sedimentary basins, or specific strata are usually studied as research objects. Recording the samples position solely based on spatial coordinates cannot meet the needs of scientific research. While relying exclusively on high-precision latitude–longitude coordinates is insufficient for rigorous spatial analyses, such attributes are nevertheless the only consistently available resource in most databases and thus remain the primary handle for sample positioning. The decimal format of latitude and longitude coordinates has been considered the most suitable recording format in the previous zircon databases (Puetz et al., 2021, 2024a, 2024b). However, a considerable number of research papers report coordinates in the DMS (Degree-Minute-Second) format. To expedite the standardization of these diverse DMS notations into a decimal format, we have crafted and implemented a Python code snippet, as detailed in Supplement 3. Another challenge arises from the absence of coordinate reports in some papers. Traditionally, papers lacking specific coordinates have been excluded from databases. However, directly exclusion could exacerbate the spatio-temporal bias. **To enhance the data richness, a spatial coordinate estimation method was applied during the database construction process.** This method, based on a plane graph and implemented in Python, swiftly estimates coordinates for articles missing these details while striving to maintain accuracy (Supplement 3).

Given the significance of detrital zircons in geological research, the strata information schema within our database has been designed to encapsulate a wide array of sedimentary data. It documents the strata age according to the period-epoch-stage stratigraphic system, as well as the maximum, estimated, and minimum depositional ages. Further details regarding the stratigraphic data points are outlined in Table 4.

Although maximizing the utilization of research papers can mitigate spatial bias to a certain extent, the spatial-strata information visualized in both the U-Pb and Lu-Hf (Figure 2-4) datasets continues to exhibit significant spatial skew. A majority of the records are concentrated in East Asia, with a particular focus on China.

Despite this concentration, all indicators suggest a substantial global representation within our datasets. The visualization tools employed highlight the areas of high research activity while also underscoring the need for further research in underrepresented regions to achieve a more balanced global perspective.

### 3.3 U-Pb isotopes database

The geochronological records include full analytical-method metadata: the technique used (e.g., LA-ICP-MS, SHRIMP or ID-TIMS), the analytical institution's spot location (rim vs core), and the spot diameter. For the chronological data, the isotopic ratios  $^{206}\text{Pb}/^{238}\text{U}$ ,  $^{207}\text{Pb}/^{235}\text{U}$ ,  $^{207}\text{Pb}/^{206}\text{Pb}$ , and  $^{208}\text{Pb}/^{232}\text{Th}$  were recorded with corresponding  $1\sigma$  uncertainties. A limited number of papers have reported uncertainties at the  $2\sigma$  level. Where a preferred age was not explicitly reported by the original authors ( $\leq 0.5\%$  of records), OneDZ estimated the most reliable date using  $1\sigma/2\sigma$  uncertainty and the 1200 Ma/1600 Ma thresholds

of Gehrels et al. (2008). These rare “estimated ages” are flagged as EstAge = 1 so that users can readily distinguish them from author-specified values.

230 Furthermore, the database also archives the discordance ratio, concentrations of U, Th, and Pb, as well as the Th/U ratios, providing a comprehensive set of parameters for geochronological analysis.

### 3.4 Lu-Hf isotopes database

The Lu-Hf isotopic data within OneDZ are fundamentally anchored in U-Pb chronological results. Alongside the age determinations, we have meticulously documented the basic analytical results, including the  $^{176}\text{Yb}/^{177}\text{Hf}$ ,  $^{176}\text{Lu}/^{177}\text{Hf}$ , and  
235  $^{176}\text{Hf}/^{177}\text{Hf}$  isotopic ratios, each accompanied by their corresponding  $2\sigma$  uncertainties, which reflect the precision of the measurements.

In addition to the raw isotopic data, OneDZ encompassed several calculated parameters derived from these ratios. These include the calculated ratios of the hafnium isotope composition  $\epsilon\text{Hf}(t)$  with their respective  $2\sigma$  uncertainties, and the model ages TDM1 (Ma) and TDM2 (Ma). These calculated results provide further insights into the isotopic evolution and the crustal  
240 residence history of the samples analyzed.

## 4 Data characteristics

### 4.1 Rock types statistics

Clastic sediments are vital geological archives to offer deep insights into the sedimentary provenance and evolutionary history of the continental crust (Taylor, 1985). In preparation for studies on sedimentary provenance and related geological inquiries,  
245 the OneDZ database collected petrological contexts from original articles. OneDZ categorized rock types into the widely accepted hierarchical granularity system, with Class-1 encompassing clastic, meta-clastic, and pyroclastic rocks. These categories reflect the diverse origins of sediments. Specifically, Class-2 and Class-3 types provide a more nuanced classification based on grain size, which is crucial for understanding sedimentary processes and environments. Class-2 further subdivided the rocks, serving as a supplement to the macroscopic rock classification of Class-1. Class-3 adopts the particle  
250 size classification scheme for detrital sedimentary rocks proposed by Udden (1918), Wentworth (1922), and Krumbein (1938), and provides the most detailed classification of rock types. In the U-Pb datasets, Class-1 rock types are predominantly clastic (50%, Figure 5a). Meta-clastics are the second lithological source (36.3%, Figure 5a). Pyroclastic takes up a little ratio (13.8%, Figure 5a). For Class-2 rock types, the major component is sandstone (53.6%, Figure 5b). The breccia, shale, mudstone equally allocated the remaining proportion (13%, 15.7%, 17.8%, Figure 6b). For Class-3 rock types, the distribution of particle sizes  
255 is quite uniform, with the proportions ranging from 1.8% to 16.0%. Specifically, fine sand and very coarse sand are the predominant types, accounting for 16.0% and 12.5% of the total, respectively (Figure 5c).

In the Lu-Hf datasets, relatively little rock records were provided in the original article. Clastic rock types contributed the most data to the dataset, comprising 38.9% of the total (Figure 5d). Meta-clastic offered 28.4% of the data and pyroclastic provided

32.7% of the data (Figure 5d). For Class-2 rock types, the major component is sandstone (66.8%, Figure 5e). The breccia, shale, mudstone equally allocated the remaining proportion (9.2%, 11.3%, 12.8%, Figure 5e). In the Class-3 rock types, the distribution of grain sizes is dominated by fine sand, which accounts for 37.5% of the total, making it the most prevalent grain size (Figure 5f). Very coarse sand is also a significant component, comprising 9.3% of the dataset. Other grain sizes contribute with percentages ranging from 1.4% to 8.0%, indicating a relatively balanced but varied composition across the different grain sizes.

## 265 4.2 Data uncertainty

The data uncertainty in the OneDZ database is stemmed from methodological errors, dating uncertainties, and potential biases associated with analytical instruments. Methodological errors are primarily attributed to variations in decay constants and half-lives among different isotopic systems. Dating uncertainties and potential biases have more to do with data processing. Figure 6 illustrates the relationships between isotopic ratios, calculated ages, and their corresponding  $2\sigma$  uncertainties. The  $^{206}\text{Pb}/^{238}\text{U}$  isotopic system adheres to a first-order linear regression model (Figure 6a), demonstrating a relatively consistent uncertainty across a wide range of ages. However, for samples with depositional ages exceed approximately 2000 Ma, the  $2\sigma$  uncertainty of ages increases to around 300 Ma. This trend suggests that approximately 67% of samples may be associated with a temporal uncertainty of approximately 600 Ma. In contrast, the  $^{207}\text{Pb}/^{235}\text{U}$  and  $^{207}\text{Pb}/^{206}\text{Pb}$  isotopic systems are characterized by second-order polynomial regressions (Figure 6b-c). The complex regression models suggest a  $2\sigma$  uncertainty of 500 Ma emerging at around 3000 Ma in  $^{207}\text{Pb}/^{235}\text{U}$  and  $^{207}\text{Pb}/^{206}\text{Pb}$  isotopic systems. The age uncertainty becomes significantly pronounced when analyzing samples over 3000 Ma. Because the uncertainties in  $^{207}\text{Pb}/^{235}\text{U}$  and  $^{207}\text{Pb}/^{206}\text{Pb}$  isotopic systems accumulate with time and become significant only in very old samples, these systems are best suited for dating ancient rocks. (between 1000 and 3000 Ma). The relatively low uncertainties suggest  $^{207}\text{Pb}/^{235}\text{U}$  and  $^{207}\text{Pb}/^{206}\text{Pb}$  isotopic systems are particularly valuable for studying the early history of the earth's crust.

280 In addition to the intrinsic variability of isotopic systems, dating uncertainty in OneDZ database is significantly influenced by the selection of the best age. Figure 7 provides a visual representation of the discrepancies between calculated isotope ages and the best ages selected from the raw data extracted directly from published papers. Dating uncertainties grow with increasing best ages across all isotopic systems (Figure 7a-c). To address this, we employed advanced statistical techniques, including Monte-Carlo resampling (Figure 7d-f) and Bootstrap resampling (Figure 7g-i), coupled with locally weighted scatter plot smoothing (LOWESS) to estimate and visualize the dating uncertainties. The LOWESS trend lines indicate that potential thresholds of uncertainty may lie around 1000 Ma and 3000 Ma. Samples younger than 1000 Ma exhibit minimal bias, suggesting that the choice of isotopic system and the application of resampling methods have a limited impact on data uncertainty. However, isotopic uncertainties compound with time and reach  $\sim 500$  Ma by ages of 3000 Ma, limiting reliable dating to still older samples. While commonly employed strategies such as filtering samples based on acceptable  $2\sigma$  uncertainty or utilizing resampling techniques aim to mitigate the adverse effects of selecting the best age, the analysis presented in Figure 7 suggests that filtering alone does not significantly reduce uncertainties associated with the best age. Notably, in all isotopic

systems, filtered results often reveal substantial gaps in the best age, indicating that the filtering process may not be sufficient to address the underlying uncertainties. The resampling methods, however, demonstrate a capacity to alleviate these gaps, particularly in the  $^{207}\text{Pb}/^{206}\text{Pb}$  isotopic system, where they prove effective in reducing the best age discrepancies (Figure 7i).

295 For analytical techniques, the LA-ICP-MS has become the preferred method for sedimentary research due to its efficiency in yielding geochronological data. Figure 8 illustrates the variation in discordance ratios over time for different analytical instruments. The SHRIMP method, known for its precision, demonstrates a consistently low discordance ratio (Figure 8a). Remarkably, even samples with elevated age uncertainties maintain discordance ratios below 0.5%, indicating SHRIMP's reliability in dating. LA-ICP-MS displays an increase in age uncertainty for samples exceeding 1000 Ma but maintains a

300 discordance ratio below 0.5% for these samples (Figure 8b). However, LA-ICP-MS exhibits a notable disadvantage for samples with low age uncertainties dating from approximately 800 to 1200 Ma, where the discordance ratio can be relatively high, occasionally exceeding 1%. This underscores the need for particularly careful data interpretation in these specific age ranges. The ID-TIMS method, while less commonly utilized in sediment dating, exhibits low discordance ratios (Figure 8c). This suggests that ID-TIMS, despite its limitations, offers robust results for the most precision of dating requirements. Samples

305 analyzed by SIMS appear to exhibit a potential linear relationship between age uncertainty and discordance ratio (Figure 8d). The original uncertainty associated with the Lu-Hf dataset predominantly pertains to the analytical outcomes obtained from isotopic measurements. Across all geological periods, the  $2\sigma$  errors for the isotopic ratios  $^{176}\text{Hf}/^{177}\text{Hf}$ ,  $^{176}\text{Lu}/^{177}\text{Hf}$ , and  $^{176}\text{Yb}/^{177}\text{Hf}$  typically fluctuate around  $2 \times 10^{-5}$ , as depicted in Figure S7 in the s Supplement. The measurement ranges for these three isotopes are approximately  $2 \times 10^{-2}$ , indicating a high level of precision in the analytical process (Figure S7 in the

310 Supplement). The uncertainties for the Lu-Hf isotopic system are notably an order of magnitude lower than the analytical results, suggesting that the system is inherently more precise than the measurements themselves. This stability in Lu-Hf uncertainties is maintained even at high resolutions, highlighting the robustness of the dataset in providing reliable isotopic age estimates. Other uncertainty in Lu-Hf datasets are the  $\epsilon\text{Hf}(0)$  and  $\epsilon\text{Hf}(t)$  errors (Figure S8 in the Supplement). The high-quality isotopic results obtained from the Lu-Hf dataset contribute to the stable and low  $2\sigma$  errors observed in both  $\epsilon\text{Hf}(0)$  and

315  $\epsilon\text{Hf}(t)$ , as depicted in Figure S8 in the Supplement. These parameters reflect the hafnium isotope composition at the time of zircon crystallization ( $\epsilon\text{Hf}(0)$ ) and at a specific time in the past ( $\epsilon\text{Hf}(t)$ ), exhibiting a consistency in error magnitude that underscores the reliability of the dataset. Similar to the isotopic uncertainty observed in the Lu-Hf system, the uncertainties associated with  $\epsilon\text{Hf}(0)$  and  $\epsilon\text{Hf}(t)$  are considerably larger than their corresponding  $2\sigma$  errors. This discrepancy highlights the precision of the isotopic measurements relative to the calculated uncertainties of the hafnium isotope ratios. The stability of

320 the error over the timescale is particularly noteworthy, suggesting that  $\epsilon\text{Hf}(0)$  and  $\epsilon\text{Hf}(t)$  values are independent and robust indicators of the isotopic evolution of the samples. This temporal stability further reinforces the reliability of these parameters in geochronological and geochemical analyses.

### 4.3 Spatial and temporal distributions of samples

Spatial distribution biases are evident within the OneDZ database (Figure 3-4). To delve into the effects of biased distributions, the U-Pb age data was segmented according to geological time sequences and visualized (Figure 9). Temporal slices reveal that the Qinghai-Tibet Plateau, Alps, Cordillera and Andes mountains are the main sampling areas in the Cenozoic (Figure 9a-c). In the Mesozoic, the main sampling regions are similar to areas from the Cenozoic (Figure 9d-f). In the Paleozoic, East Asia is obviously over-sampled relative to other regions (Figure 9g-l). In the pre-Cambrian period, East Asia, Europe and Australia contributed **the majority of samples** (Figure 9m-n)

In this study, we also present the visualization of the temporal distributions of uranium, thorium, and lead concentrations in detrital zircons (Figure 10). The concentrations of these elements exhibit temporal stability, with uranium ranging from approximately 100 ppm to 300 ppm, thorium from 100 ppm to 200 ppm, and lead from 0 ppm to 200 ppm. Notably, there are differences in the estimation of temporal distributions of element concentrations when using Bootstrap and Monte Carlo methods (Figure 10). Furthermore, beyond elemental concentrations, the Th/U ratio in zircon is a crucial indicator for determining the provenance of zircon. It is widely accepted that a Th/U ratio below 0.1 suggests zircon may have experienced metamorphism and recrystallization, while a ratio above 0.4 is indicative of magmatic zircon. The resampling methods displays from all time spans, that the Th/U is larger than 0.4, indicates magmatic zircon dominants the detrital zircon (Figure S9 in the Supplement).

For Lu-Hf isotopes, the  $^{176}\text{Hf}/^{177}\text{Hf}$  isotope decreases with the  $^{176}\text{Lu}/^{177}\text{Hf}$  and  $^{176}\text{Yb}/^{177}\text{Hf}$  isotopes showing periodic fluctuations (Figure 11). The  $\varepsilon\text{Hf}(0)$  displays a continuous decline, while the  $\varepsilon\text{Hf}(t)$  periodically fluctuates (Figure 12).

## 5 Discussion

### 5.1 Evaluate the paleo-spatial reconstruction

Despite the OneDZ database compiles comprehensive information about detrital zircon data, obvious oversampling bias exists in regions such East Asia due to disparities in research intensity and focus. This oversampling creates an imbalance and potentially leads to overrepresentations of regional samples.

For instance, the spatial analysis of the global zircon oxygen isotope record has shown that the temporal anomalies in zircon oxygen isotopes were predominantly attributed to regional samples' imbalance (Sundell et al., 2024). To address the issue of regional disparities in global zircon data analysis, Puetz et al. (2024) introduced a method to assess global representativeness. This method involves overlaying a grid across the Earth's surface, dividing each item into discrete cells. The degree of global representativeness is then calculated by determining the ratio of the number of cells containing zircon data to the total number of cells in the grid. This approach allows for a quantitative measure of how well the zircon data cover different geographical regions. However, this evaluation approach is predicated on the present-day distribution of land and sea. In geological time scales, current geographical pattern does not accurately reflect the samples' spatial positions during the depositing period. To

enhance the analysis of the spatiotemporal representativeness, we undertook a reconstruction based on the spatial distribution of detrital zircon U-Pb data. Utilizing tools such as pyGplate (Müller et al., 2018; Mather et al., 2024) and in situ block reconstruction methods (Jian et al., 2022), samples were reconstructed following the geohistorical spatial distribution. As shown in Figure 13, the scatter plot of reconstructed data shows that OneDZ covers almost all major continents in various periods of Earth's evolution. However, the spatial kernel density map in Figure 13 re-evaluated the global representativeness of the data. In fact, as we delve into more ancient geological periods, the sampling locations tend to cluster around one or two ancient tectonic plates. This pattern is due to the fact that older zircon grains, which have undergone multiple episodes of sedimentary recycling (where the age difference between the zircon and the time of deposition exceeds 150 Ma), have been subject to significant transport and thus may not accurately represent the original paleogeographic context. Consequently, after approximately accounting for the effects of sedimentary recycling, the data from these ancient times are predominantly sourced from a limited number of locations, lacking global representativeness. Therefore, the evaluation of results based on OneDZ, indicate that the global scope of zircon big data research needs further assessment.

Following the new paleo-spatial reconstruction evaluating methods, the temporal globality of OneDZ detrital zircon U-Pb data was visualized in Figure 14. Instead of no significant variation with 2° and 4° grids (Figures 14a-b), the visualizations in Figures 14c-e demonstrate that the U-Pb data has achieved spatial coverage across paleo-continent. A notable rise (14%) in paleo-spatial reconstruction and valuable stability were observed when the grid size was enlarged from 6° to 10°. As the grid size increases, the spatial resolution of globality gradually decreases, resulting in a continuous increase in the calculated global representative values. The global representative value loses a significant amount of spatial detail when calculated at an excessively large scale. Similarly, small-scale grid calculation results in computational bias towards local detail information, leading to underestimation of the globality. After considering both local and global information, 6° is deemed suitable for evaluating the global representativeness of U-Pb data in the OneDZ database.

Figures 14c-e also show periodic peaks in globality coincides with specific geological eras. This phenomenon might be correlated with the heightened research interest in these periods. Samples from these periods are more likely to stimulate scientific inquiry due to the dynamic geological processes occurring at those times. Despite the large volume of data in OneDZ, the calculated paleo-spatial reconstruction does not fully represent global features for most geological times, as the reconstruction does not account for 100 percent of the spatial details. For example, the sample distribution reconstructed for 250 Ma (as shown in Figure 13g-h) appears to have global coverage. However, the calculated paleo-spatial reconstruction for this period only accounts for approximately 30% to 60% of the actual spatial distribution. Consequently, we recommend that regional data be handled with greater caution when interpreting global geological events. It is particularly important to employ spatial kernel density evaluation methods to ensure a more accurate representation of the data.

## 5.2 Compare the resampling methods

The temporal evolution of zircon U-Pb data is often analysed through big data methods and plays a crucial role in understanding the development of orogenic belts and crustal thickness. Big data methods with zircon U-Pb offer insights into Earth system

evolution based on anomalies in time series data. Usually, the fluctuations in the curve are explained as the evolution of the Earth system. Not only is there a risk of data not being globally representative, but the zircon U-Pb curves obtained from big data analysis may also be statistically biased due to inconsistent data volumes. Some resampling statistical tools like Bootstrap and Monte-Carlo methods are applied in zircon big data analysis (Keller & Schone, 2012; Yang et al., 2024; Yang et al., 2025). These methods have usually been assumed to be effective in previous studies. However, these resampling methods have not been systematically tested. The zircon U-Pb data in OneDZ, as the world's largest multidimensional imbalanced spatiotemporal dataset, provides a data foundation for comparing the effects when applying different resampling methods.

Firstly, we selected the best age data from zircon U-Pb data for time resampling experiments. In addition to comparing Bootstrap and Monte-Carlo resampling methods, we assessed the impact of data sparsity using the  $2\sigma$  error to identify potential outliers and quantify the uncertainty. The experiment focuses on the sparsity of samples generated within the time range of zircon U-Pb ages exceeding 2500 Ma, with a threshold of 400 Ma. After time resampling using Monte-Carlo (Figure 8d-f) and Bootstrap methods (Figure 8g-i), the overall trend of zircon best age data is consistent. Even on time series after 2500 Ma, there was no significant difference in the characterization of evolutionary trends between the two resampling methods.

However, there is a significant difference between the two methods in characterizing the details of a time series. In the  $^{206}\text{Pb}/^{238}\text{U}$  isotope system, four periodic fluctuations were observed in the Monte-Carlo resampling results over the time period of 0-1000 Ma (Figure 8d). The Bootstrap method only shows a slight increase around 500 Ma on the same time scale (Figure 8g). The rest of the time scales show a slow increase. In the  $^{207}\text{Pb}/^{235}\text{U}$  isotope system, the Monte-Carlo resampling results showed four small amplitude periodic fluctuations in the 0-2000Ma time period under a generally slow rising background (Figure 8e). The Bootstrap method showed a significant decrease around 1500 Ma on the same time scale (Figure 8h). In the  $^{207}\text{Pb}/^{206}\text{U}$  isotope system, the Monte-Carlo resampling results showed a significant decrease around 1500Ma (Figure 8f). In contrast, the Bootstrap method shows a periodic decrease (Figure 8i), which differs from the more substantial decrease observed with other methods. Although Figure 8 overall depicts the magnitude of age error over time in different systems and does not have practical geological significance, the significant differences in the time curves after resampling using Monte-Carlo and Bootstrap methods indicate the need for caution in interpreting data after applying resampling methods. Furthermore, we compared the results of time resampling methods for zircon U-Pb and Lu-Hf system time series data in OneDZ. In the analysis of zircon U-Pb data, the Bootstrap method demonstrates greater consistency over time (Figure 10a-c), meaning that the results obtained using this method exhibit less variation across different time periods compared to other methods. The Monte-Carlo method is more sensitive to local data fluctuations than the Bootstrap method (Figure S7 in the Supplement). The Monte Carlo method also shows significant oscillations on relatively sparse  $\varepsilon\text{Hf}(0)$  and  $\varepsilon\text{Hf}(t)$  and corresponding errors data (Figure S7a-c in the Supplement). The difference between Bootstrap and Monte Carlo methods will also disappear as the amount of data increases. In the  $^{176}\text{Yb}/^{177}\text{Hf}$   $2\sigma$  error time series, due to the significant increase in data volume, the significant difference in the results of resampling methods is relatively little (Figure S8 in the Supplement). The above experimental time series data density statistics show that different resampling methods are actually controlled by data density and the areas where significant oscillations occur in the Monte-Carlo method coincide with areas with high data density (Figure 10-11, S7-S7).

Since standard Monte-Carlo time-resampling assumes that the underlying process is stationary (often approximated by normality or local uniformity, by Rubinstein & Krose (2016)), it can yield biased estimates when the data density evolves sharply within the chosen window, which is a situation common in high-frequency zircon datasets. Consequently, more flexible, density-aware resampling strategies are preferred for zircon big-data analysis.

425 **Spatial over-sampling introduces another potential bias that has gained attention in the field (Keller et al., 2018).** Addressing this issue often involves spatial resampling methods, which were employed in this research using the OneDZ database. Initially, Monte Carlo spatial resampling was used to assess the frequency at which samples are selected (Keller et al., 2018). Ideally, a balanced spatial sampling should achieve equal total sampling frequencies across regions, increasing the likelihood of sampling from underrepresented areas. Our findings suggest that direct application of the Monte-Carlo method does not mitigate  
430 sampling bias. Samples from East Asia, particularly China, remain overrepresented due to the large volume of available data from this region, skewing the overall data distribution and leaving other regions sparsely represented, similar to the observed sample sparsity in the temporal domain (Figure 15a). To counteract the hypothesis, we explored data augmentation methods to generate new data points in under-sampled regions. This study introduces the Synthetic Minority Over-sampling Technique (SMOTE, Chawla et al., 2002) to create synthetic data points from regions other than China while preserving the same data  
435 features. Applying SMOTE led to a significant increase in resampling frequency in these regions (Figure 15b). Inspired by grid-based methods, we also pre-processed the data by averaging the U-Pb age signals before applying SMOTE. This novel approach enhanced resampling frequency in previously under-sampled regions, resulting in the sampling differences in different regions to significantly reduce (Figure 15c). Direct spatial resampling methods may not adequately resolve spatial imbalances. However, combining these methods with data enhancement techniques and grid-based approaches can  
440 significantly mitigate spatial biases.

### 5.3 Implications for database construction and future developments

**The compilation of the OneDZ dataset, which employs a crowdfunding approach, has the potential to significantly broaden data coverage (comparison with other databases can be seen in Figure S10). However, crowdfunding introduces challenges, such as inconsistencies in data formatting and the risk of human errors. To address these issues, a series of Python scripts for automated data cleaning and inspection were developed. These scripts have successfully replaced labor-intensive manual inspections, reducing both labor costs and the potential for data errors. From the OneDZ compilation process, crowdfunding combined with automatic data cleaning by Python code snippets is feasible and greatly improved the efficiency of dataset assembly.**

Moreover, AI tools have played a pivotal role in the data collection and extraction process. Unlike traditional web crawlers,  
450 which can pose privacy risks, AI models can predict whether an article may contain the required database features based on publicly available text information, such as titles and abstracts. The integration of AI models into the database construction process eliminates the need for manual screening of potential articles, significantly improving efficiency. Additionally, computer vision tools like DeepShovel are crucial, as a considerable amount of article data is stored in PDF image files in the

form of tables. Manual reading and data storage are not sustainable approaches for handling such large volumes of data.

455 Computer vision-based AI models show great promise in reducing labor costs and increasing efficiency.

As of the current submission, the number of U-Pb entries in this dataset has increased from approximately 1.8 million in the initial preprint to 2.5 million, suggesting that eliminating manual data organization and template-based input substantially encourages community contributions. However, this rapid expansion introduces notable challenges in data quality and source relevance. Although agents demonstrate high efficiency in extracting data from heterogeneous source files, including manuscript PDFs and supplementary materials in Word, Excel, or PDF formats (Figure S11 in the Supplement), the automated workflow occasionally yields extraction errors, making manual verification indispensable (Supplement 5). Nevertheless, given the time-intensive nature of meticulous checking, verified records remain at 1.5 million, considerably fewer than the total collected entries, which necessitates clearly distinguishing checked and unchecked records upon publication and explicitly cautioning users regarding unchecked data. Moreover, without template constraints, users may upload arbitrary PDF files unrelated to detrital zircon studies, prompting the deployment of an independent classification agent to assess file relevance, whereby over five thousand irrelevant publications, such as those concerning magmatic or granitic zircons and general geochemical articles, were identified and filtered during our experiments (Figure S12 in the Supplement). Therefore, integrating large language models into database construction does not inherently reduce manual labor but rather demands additional human oversight to ensure data quality and reliability.

460

465

470 Future work will focus on refining the automated extraction pipelines and expanding community-contributed data coverage, rather than on developing custom web portals.

### Data availability

The complete OneDZ dataset (Li et al., 2026) is freely available as flat CSV files via Zenodo at <https://doi.org/10.5281/zenodo.19690702>. The release includes 22 sequential CSV parts for the full U-Pb compilation (2,550,738 records, 64 standardised columns), 3 parts for the Lu-Hf compilation (297,527 records, 33 standardised columns), and 14 parts for the expert-verified U-Pb subset (1,414,062 records). All files use UTF-8 encoding with comma delimiters and a single header row, enabling direct ingestion without database infrastructure. For users who prefer to import the data into relational database environments, an archived SQL dump of the same dataset is separately available at <https://doi.org/10.5281/zenodo.17407937> (this is not required to access or use the data). Python code snippets for data cleaning and harmonisation used in this research are accessible via <https://github.com/KeranLi/Global-Detrital-Zircon>. For researchers who wish to replicate the LLM-driven extraction workflow, the original data sources and agent outputs are accessible at <https://doi.org/10.5281/zenodo.19691004>, and the corresponding code snippets are available at [https://github.com/KeranLi/onedz\\_llm\\_coding](https://github.com/KeranLi/onedz_llm_coding).

475

480

## Conclusion

485 In this study, we introduce a ground-breaking global detrital zircon U-Th-Pb geochronology and Lu-Hf isotope database, which serves as a critical resource for advancing Earth science research. This database includes 1925687 U-Pb and 275971 Lu-Hf records, offering a broad sampling range from global detrital rocks. The database offers an extensive and diverse collection of data, including various types of stratigraphic information, a broad range of sedimentary ages, comprehensive isotope geochemical datasets, and data from multiple analytical techniques such as LA-ICP-MS, SHRIMP, SIMS, and TIMS.

490 Based on this database, we have characterized the uncertainties associated with zircon dating, compared the efficacy of different analytical techniques, proposed an evaluation method that assesses the deep-time global coverage of the data and discussed the challenges and potential solutions related to spatiotemporal sampling methods. Although the data is globally sourced, variations in spatial and temporal distribution can affect its global representativeness. Therefore, when conducting big data analyses on spatial or temporal distributions, reconstructing data's paleo-points is necessary. In imbalanced

495 spatiotemporal data resampling methods, Bootstrap methods and SMOTE data augmentation methods may be more suitable. **The development of OneDZ demonstrates that leveraging crowdfunding, large language models, and automated code cleaning processes is essential for the rapid assembly of a comprehensive geoscience dataset. By distributing the final product as standardised, flat CSV files via Zenodo, we ensure that the data are immediately accessible to any researcher without requiring proprietary database software or custom web interfaces. Furthermore, integrating AI tools and scripting workflows enhances**

500 **both the efficiency of data extraction and the reliability of harmonised outputs, making large-scale compilations more reproducible and transparent.**

## Author contributions

KL, XH, RC, JH, WX, YP, AM, HH, QG, WY, LH, LQ, GC, GS, SZ, TD, KL, JS and BG compiled the data. KL, RC and WX merged the data, formatted the data, performed the analyses, standardized the reference materials, organized the database,

505 managed the publication of the database in the Zenodo repository, and drafted and revised the manuscript. KL designed the code snippets. TL and CF developed the web platform. HX initiated and supported the data compilation.

## Competing interests

The contact author has declared that none of the authors has any competing interests.

## Acknowledgements

510 The authors would thank Dr. S. J. Puetz and Wu Y. and their colleagues for establishing 1.2 million detrital zircon database and Chinese zircon database. We thank the sedimentary group members in IUGS Deep-time Digital Earth (DDE) Big Science

Program for their assistance in collecting and cleaning detrital zircon data in China. The related sub database research has been published in the special issue of the *Geoscience Data Journal* in 2024. This work was financially supported by the National Natural Science Foundation of China (42142004). This paper contributes to the IUGS “Deep-time Digital Earth” Big Science Program. This paper got support from high performance computing center, Nanjing University in reconstructing the paleo-locations of records.

## References

- Chai, R., Yang, J., Deng, T., and Hu, X.: A detrital zircon dataset for the eastern Central China Orogenic Belt (East Qinling, Dabie and Sulu orogens), *Geoscience Data Journal*, 11, 4, 562-572, <https://doi.org/10.12297/dpr.dde.202212.3>, 2023.
- 520 Chawla, N. V., Bowyer, K. W., Hall, L. O., and Kegelmeyer, W. P.: SMOTE: synthetic minority over-sampling technique, *Journal of artificial intelligence research*, 16, 321-357, <https://doi.org/10.1613/jair.953>, 2002.
- Chen, G., Li, C., Shi, Y., and Zha, K.: A synthesis of available detrital zircon data from Turkey, Cyprus and Greek peninsula, *Geoscience Data Journal*, 11, 2, 137-147, <https://doi.org/10.1002/gdj3.216>, 2023a.
- Chen, X., Wang, P., Xie, H., Zhu, L., Liao, X., and Kong, X.: Detrital zircon U-Pb ages and Hf isotope analyses of modern and Quaternary sediments in China: A new dataset with preliminary analysis, *Geoscience Data Journal*, 11, 4, 374-384, 525 <https://doi.org/10.1002/gdj3.193>, 2023b.
- Cheng, Q.: Non-linear theory and power-law models for information integration and mineral resources quantitative assessments, *Mathematical Geosciences*, 40, 503–532, 10.1007/s11004-008-9172-6, 2008.
- Chu, X., Ilyas, I. F., Krishnan, S., and Wang, J.: Data cleaning: Overview and emerging challenges, in: Proceedings of the 530 2016 international conference on management of data, pp. 2201–2206, <https://doi.org/10.1145/2882903.2912574>, 2016.
- Claesson, S., Vetrin, V., Bayanova, T., and Downes, H.: U–Pb zircon ages from a Devonian carbonatite dyke, Kola peninsula, Russia: a record of geological evolution from the Archaean to the Palaeozoic, *Lithos*, 51, 95-108, [https://doi.org/10.1016/S0024-4937\(99\)00076-6](https://doi.org/10.1016/S0024-4937(99)00076-6), 2000.
- Deng, C., Jia, Y., Xu, H., Zhang, C., Tang, J., Fu, L., Zhang, W., Zhang, H., Wang, X., and Zhou, C.: GAKG: A multimodal geoscience academic knowledge graph, in: Proceedings of the 30th ACM International Conference on Information & Knowledge Management, 535 pp. 4445–4454, <https://doi.org/10.1145/3459637.3482003>, 2021.
- Dong, Y., Zuo, P., Xiao, Z., Zhao, Y., Zheng, D., Sun, F., and Li, Y.: A database of detrital zircon U-Pb ages in the North China Craton from the Paleoproterozoic to the early Palaeozoic, *Geoscience Data Journal*, 11, 4, 365-373, <https://doi.org/10.1002/gdj3.192>, 2023.
- 540 Gehrels, G. E., Valencia, V. A., and Ruiz, J.: Enhanced precision, accuracy, efficiency, and spatial resolution of U-Pb ages by laser-ablation–multicollector–inductively coupled plasma–mass spectrometry, *Geochemistry, Geophysics, Geosystems*, 9, 3, Q03017, doi:10.1029/2007GC001805, 2008.

- He, W., Sun, J., Dong, Y., Qian, T., Wang, T., He, L., and Qi, Y.: A synthesis of available detrital zircon data from the Qilian-Qaidam-Kunlun collage, northern Tibet, *Geoscience Data Journal*, 11, 4, 465-478, <https://doi.org/10.1002/gdj3.225>, 2023.
- 545 Hellerstein, J. M.: Quantitative data cleaning for large databases, <http://db.cs.berkeley.edu/jmh>, 2013.
- Hoskin, P. W. and Ireland, T. R.: Rare earth element chemistry of zircon and its use as a provenance indicator, *Geology*, 28, 627–630, [https://doi.org/10.1130/0091-7613\(2000\)28<627:REECOZ>2.0.CO;2](https://doi.org/10.1130/0091-7613(2000)28<627:REECOZ>2.0.CO;2), 2000.
- <https://doi.org/10.48550/arXiv.2210.02830>, 2022a.
- Hu, X.M., Xu, Y.W., Ma, X.G., Zhu, Y.Q., Ma, C., Li, C., Lü, H.R., Wang, X.B, Zhou, C.H. and Wang, C.S.: Knowledge System, Ontology, and Knowledge Graph of the Deep-Time Digital Earth (DDE): Progress and Perspective, *Journal of Earth Science*, 34, 1323–1327, <https://doi.org/10.1007/s12583-023-1930-1>, 2023.
- Huang, Y. and Hu, L.: A database of detrital zircon U–Pb ages and Lu–Hf isotope of sediments in the South China Sea, *Geoscience Data Journal*, 11, 4, 433-442, <https://doi.org/10.1002/gdj3.218>, 2023.
- Jaffey, A., Flynn, K., Glendenin, L., Bentley, W. T., and Essling, A.: Precision measurement of half-lives and specific activities of <sup>235</sup>U and <sup>238</sup>U, *Physical Review C*, 4, 1889-1906, <https://doi.org/10.1103/PhysRevC.4.1889>, 1971.
- 555
- Jian, D., Williams, S. E., Yu, S., and Zhao, G.: Quantifying the link between the detrital zircon record and tectonic settings, *Journal of Geophysical Research: Solid Earth*, 127, e2022JB024 606, <https://doi.org/10.1029/2022JB024606>, 2022.
- Keller, C. B., and Schoene, B.: Statistical geochemistry reveals disruption in secular lithospheric evolution about 2.5 Gyr ago. *Nature*, 485(7399), 490-493, <https://doi.org/10.1038/nature11024>, 2012.
- 560
- Kinny PD, Mass R. Lu–Hf and Sm–Nd isotope systems in zircon. *Reviews in Mineralogy and Geochemistry*, 53: 327-341, <https://doi.org/10.2113/0530327>, 2003.
- Krogh, T.: A low-contamination method for hydrothermal decomposition of zircon and extraction of U and Pb for isotopic age determinations, *Geochimica et cosmochimica acta*, 37, 485-494, [https://doi.org/10.1016/0016-7037\(73\)90213-5](https://doi.org/10.1016/0016-7037(73)90213-5), 1973.
- Li K., Hu X., Chai R., Yang J., Xue W., Pan Y., Li T., Fang, C., Ma, A., Huang, H., Guo, Q., Yang, W., Hu, L., Qi, L., Chen, G., Sun, G., Zhang, S., Deng, T., Li, K., Sun, J. and Gao, B.: OneDZ: A Global Detrital Zircon Database and Implications for Constructing Giant Geoscience Database, Zenodo, <https://doi.org/10.5281/zenodo.19690702>.
- 565
- Lu, B., Wu, L., Yang, L., Sun, C., Liu, W., Gan, X., Liang, S., Fu, L., Wang, X., and Zhou, C.: DataExpo: A One-Stop Dataset Service for Open Science Research, in: Companion Proceedings of the ACM Web Conference 2023, pp. 32–36, <https://doi.org/10.1145/3543873.3587305>, 2023.
- 570
- Luo, C., Qi, L., and Xia, T.: A database of detrital zircon U–Pb ages and Hf isotope of Precambrian strata in South China, *Geoscience Data Journal*, 11, 4, 385-393, <https://doi.org/10.1002/gdj3.194>, 2023.
- Martin, E. L., Barrote, V. R., and Cawood, P. A.: A resource for automated search and collation of geochemical datasets from journal supplements, *Sci. Data*, 9, 724, <https://doi.org/10.1038/s41597-022-01730-7>, 2022.
- Mather, B. R., Müller, R. D., Zahirovic, S., Cannon, J., Chin, M., Ilano, L., Wright, N. M., Alfonso, C., Williams, S., Tetley, M., et al.: Deep time spatio-temporal data analysis using pyGPlates with PlateTectonicTools and GPlately, *Geoscience Data Journal*, 11, 1,3-10, <https://doi.org/10.1002/gdj3.185>, 2024.
- 575

- McKenzie, N.R., Horton, B.K., Loomis, S.E., Stockli, D.F., Planavsky, N.J. and Lee, C.A.: Continental arc volcanism as the principal driver of icehouse-greenhouse variability, *Science* 352, 444-447, <https://www.science.org/doi/full/10.1126/science.aad5787>, 2016.
- 580 Merdith AS, Williams SE, Collins AS, Tetley MG, Mulder JA, Blades ML, Young A, Armistead SE, Cannon J, Zahirovic S and Müller RD. Extending full-plate tectonic models into deep time: Linking the Neoproterozoic and the Phanerozoic, *Earth-Science Reviews*, 214, 103477, <https://doi.org/10.1016/j.earscirev.2020.103477>, 2021.
- Müller, R. D., Cannon, J., Qin, X., Watson, R. J., Gurnis, M., Williams, S., Pfaffelmoser, T., Seton, M., Russell, S. H., and Zahirovic, S.: GPlates: Building a virtual Earth through deep time, *Geochemistry, Geophysics, Geosystems*, 19, 2243–
- 585 2261, <https://doi.org/10.1029/2018GC007584>, 2018.
- Pan, Y. and Hu, X.: A database of detrital zircon U–Pb geochronology and Hf isotopes from the Songpan–Ganzi and Western Qinling terranes, *Geoscience Data Journal*, 11, 4, 394-404, <https://doi.org/10.1002/gdj3.195>, 2023.
- Patchett PJ, Tatsumoto M. A routine high–precision method for Lu–Hf isotope geochemistry and chronology. *Contribution to Mineralogy and Petrology*, 75: 263–267, <https://doi.org/10.1007/BF01166766>, 1981.
- 590 Patchett, P. J.: Importance of the Lu–Hf isotopic system in studies of planetary chronology and chemical evolution, *Geochimica et Cosmochimica Acta*, 47, 81–91, [https://doi.org/10.1016/0016-7037\(83\)90092-3](https://doi.org/10.1016/0016-7037(83)90092-3), 1983.
- Puetz, S. J., Condie, K. C., Sundell, K., Roberts, N. M., Spencer, C. J., Boulila, S., and Cheng, Q.: The replication crisis and its relevance to Earth Science studies: Case studies and recommendations, *Geoscience Frontiers*, 15, 101821, <https://doi.org/10.1016/j.gsf.2024.101821>, 2024a.
- 595 Puetz, S. J., Spencer, C. J., and Ganade, C. E.: Analyses from a validated global U–Pb detrital zircon database: Enhanced methods for filtering discordant U–Pb zircon analyses and optimizing crystallization age estimates, *Earth-Science Reviews*, 220, 103745, <https://doi.org/10.1016/j.earscirev.2021.103745>, 2021.
- Puetz, S. J., Spencer, C. J., Condie, K. C., and Roberts, N. M.: Enhanced U–Pb detrital zircon, Lu–Hf zircon,  $\delta^{18}\text{O}$  zircon, and Sm–Nd whole rock global databases, *Scientific Data*, 11, 56, <https://doi.org/10.1038/s41597-023-02902-9>, 2024b.
- 600 Rubinstein, R. Y., Kroese, D. P.: *Simulation and the Monte Carlo method*, John Wiley & Sons, 2016.
- Scherer, E., Carsten M., and Klaus M.: Calibration of the lutetium-hafnium clock, *Science*, 293, 5530, 683-687, 2001.
- Söderlund, U., Patchett, P.J., Vervoort, J.D., Isachsen, C.E.: The  $^{176}\text{Lu}$  decay constant determined by Lu–Hf and U–Pb isotope systematics of Precambrian mafic intrusions, *Earth Planet. Sci. Lett.* 219 (3-4), 311-324, [https://doi.org/10.1016/S0012-821X\(04\)00012-3](https://doi.org/10.1016/S0012-821X(04)00012-3). 2024.
- 605 Sun, G. and Chen, J.: A database of detrital zircon U–Pb ages and Hf isotopes for the Middle East (Iranian and Arabian plates), *Geoscience Data Journal*, 11, 2, 107-117, <https://doi.org/10.1002/gdj3.187>, 2023.
- Sundell, K. E., Macdonald, F. A., and Puetz, S. J.: Does zircon geochemistry record global sediment subduction?, *Geology*, 52, 282–286, <https://doi.org/10.1130/G51817.1>, 2024.
- Taylor, S R, M. S. M.: *The continental crust: its composition and evolution*, Black well Scientific Publications, Oxford, 1-328,
- 610 <https://commons.library.stonybrook.edu/geo-articles/12/>, 1985.

- Wang, L., Huo, N., Jiang, G., Han, C., Sun, J., and Huang, H.: Detrital zircon U–Pb and Hf isotopic dataset for the Central Asian Orogenic Belt, northern China, *Geoscience Data Journal*, 11, 4, 426-432, <https://doi.org/10.1002/gdj3.214>, 2023.
- Wu, Y., Fang, X., and Ji, J.: A global zircon U–Th–Pb geochronological database, *Earth System Science Data*, 15, 5171-5181, <https://doi.org/10.5194/essd-15-5171-2023>, 2023.
- 615 Xia, T., Li, K., Hu, L., Zhao, Z., Huang, Y., Ma, Q., and Qi, L.: A database of detrital zircon geochronology ages of Cambrian to Paleogene deposits in South China, *Geoscience Data Journal*, 11, 4, 405-413, <https://doi.org/10.1002/gdj3.196>, 2023.
- Yang, W., Li, Q., Yang, J., Fang, T., and Ma, R.: Dataset of detrital zircon U–Pb ages and Hf isotopic compositions for the late Paleozoic–Mesozoic strata in the North China block, *Geoscience Data Journal*, 11, 4, 414-425, <https://doi.org/10.1002/gdj3.211>, 2023.
- 620 Yang, X., Zhang, Z., Zhou, Y., and Yang, J.: Spatio-temporal analysis of Permian-Cretaceous magmatic activities in the Tengchong block: Implications for tectono-magmatic evolution. *Geoscience Frontiers*, 15, 6, 101920, <https://doi.org/10.1016/j.gsf.2024.101920>, (2024).
- Yang, X., Zhang, Z., Zhou, Y., Yang, J., and Wang, Y.: Decoding the geological evolution of Tengchong Block: A big data analysis of zircon from the Late Permian to Early Cretaceous. *Tectonics*, 44, 6, e2025TC008882, <https://doi.org/10.1029/2025TC008882>, 2025.
- 625 Zhang, S., Hu, X., Zhang, J., Li, Q., Xu, Y., Yu, Y., and Han, L.: A database of detrital zircon U–Pb ages and Hf isotopic compositions from the Tarim, West Kunlun, Pamir, Tajik and Tianshuihai terranes, *Geoscience Data Journal*, 11, 2, 118-127, <https://doi.org/10.1002/gdj3.213>, 2023a.
- Zhang, S., Jia, Y., Xu, H., Wang, D., Li, T. J.-j., Wen, Y., Wang, X., and Zhou, C.: KnowledgeShovel: An AI-in-the-Loop Document Annotation System for Scientific Knowledge Base Construction, arXiv preprint arXiv:2210.02830,
- 630 Zhang, S., Jia, Y., Xu, H., Wen, Y., Wang, D., and Wang, X.: Deepshovel: An online collaborative platform for data extraction in geoscience literature with ai assistance, arXiv preprint arXiv:2202.10163, <https://doi.org/10.48550/arXiv.2202.10163>, 2022b.
- Zhang, S., Xu, H., Jia, Y., Wen, Y., Wang, D., Fu, L., Wang, X., and Zhou, C.: GeoDeepShovel: A platform for building scientific database from geoscience literature with AI assistance, *Geoscience Data Journal*, 10, 519-537, <https://doi.org/10.1002/gdj3.186>, 2023b.
- 635

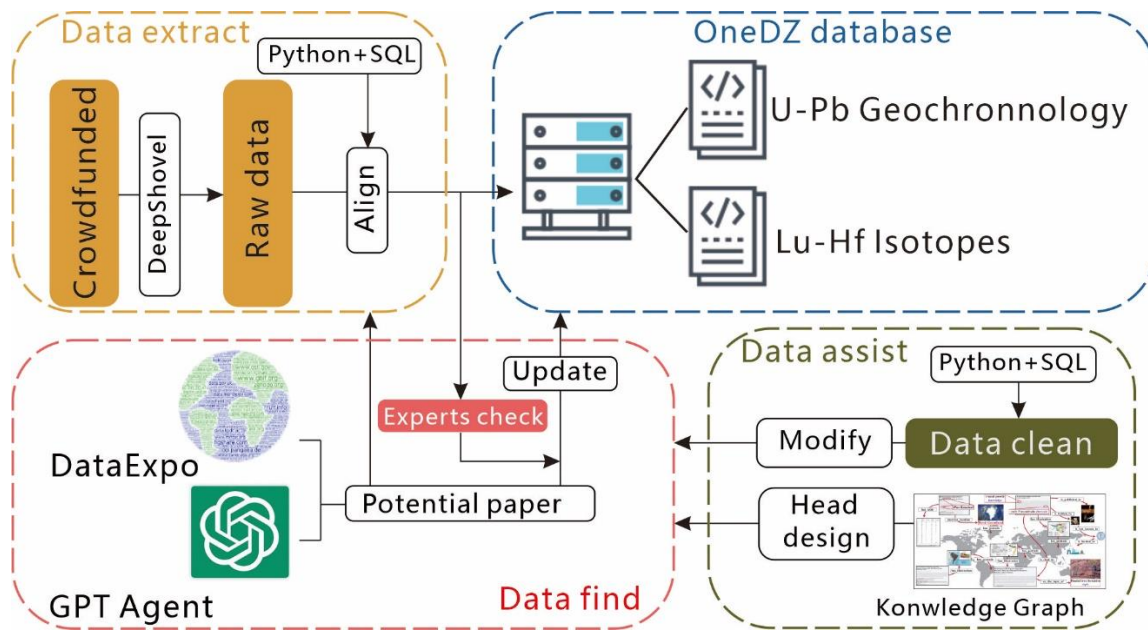
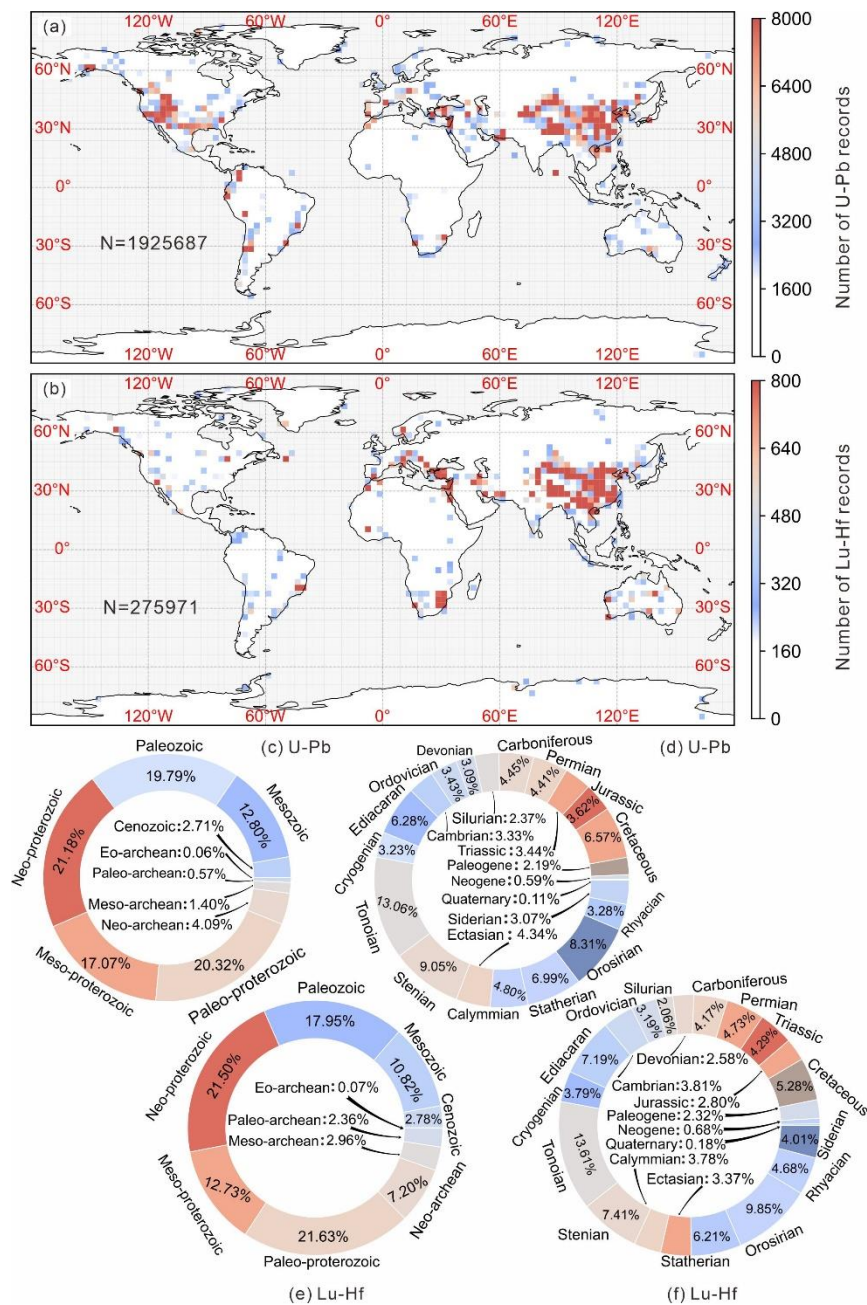
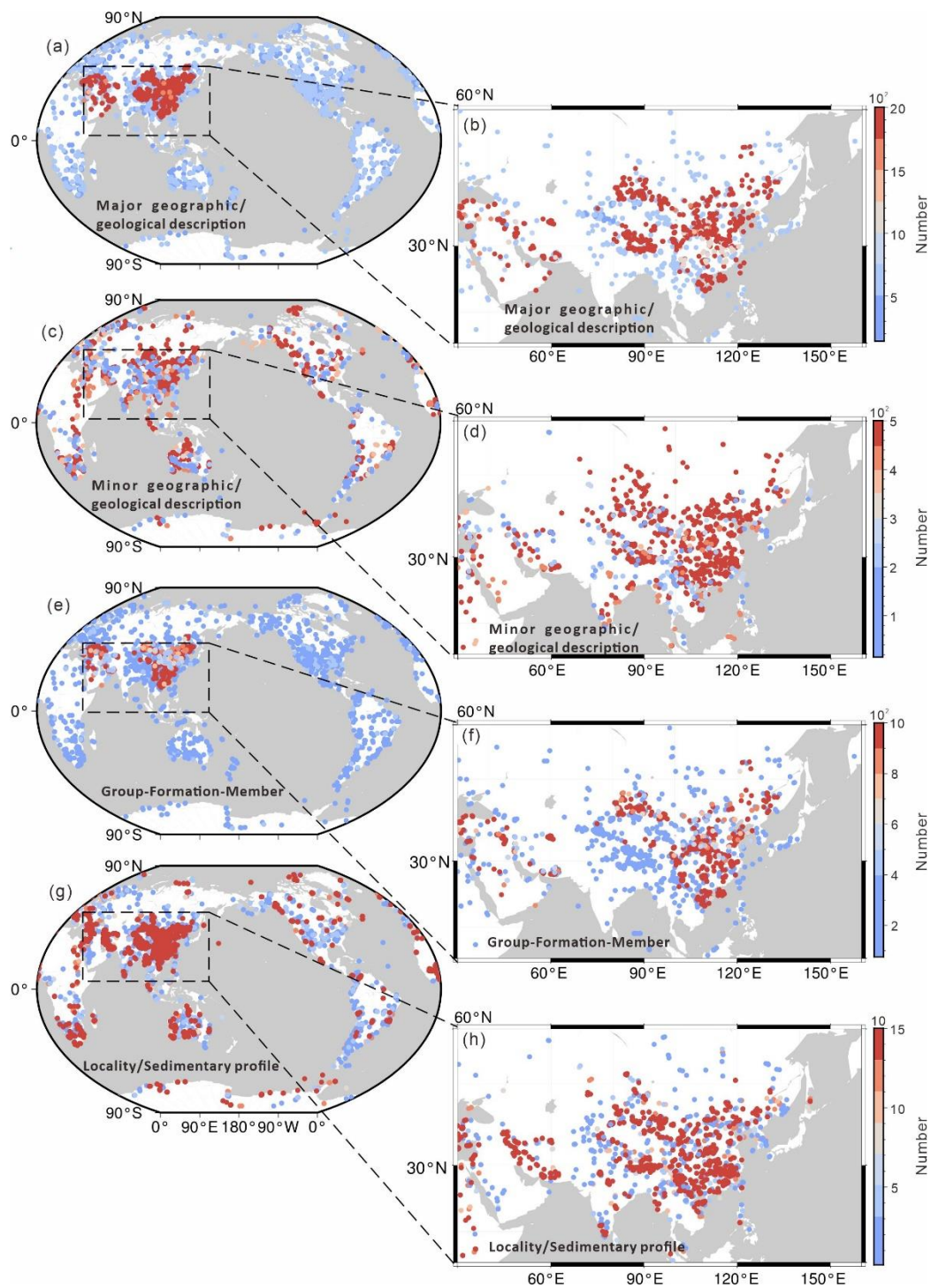


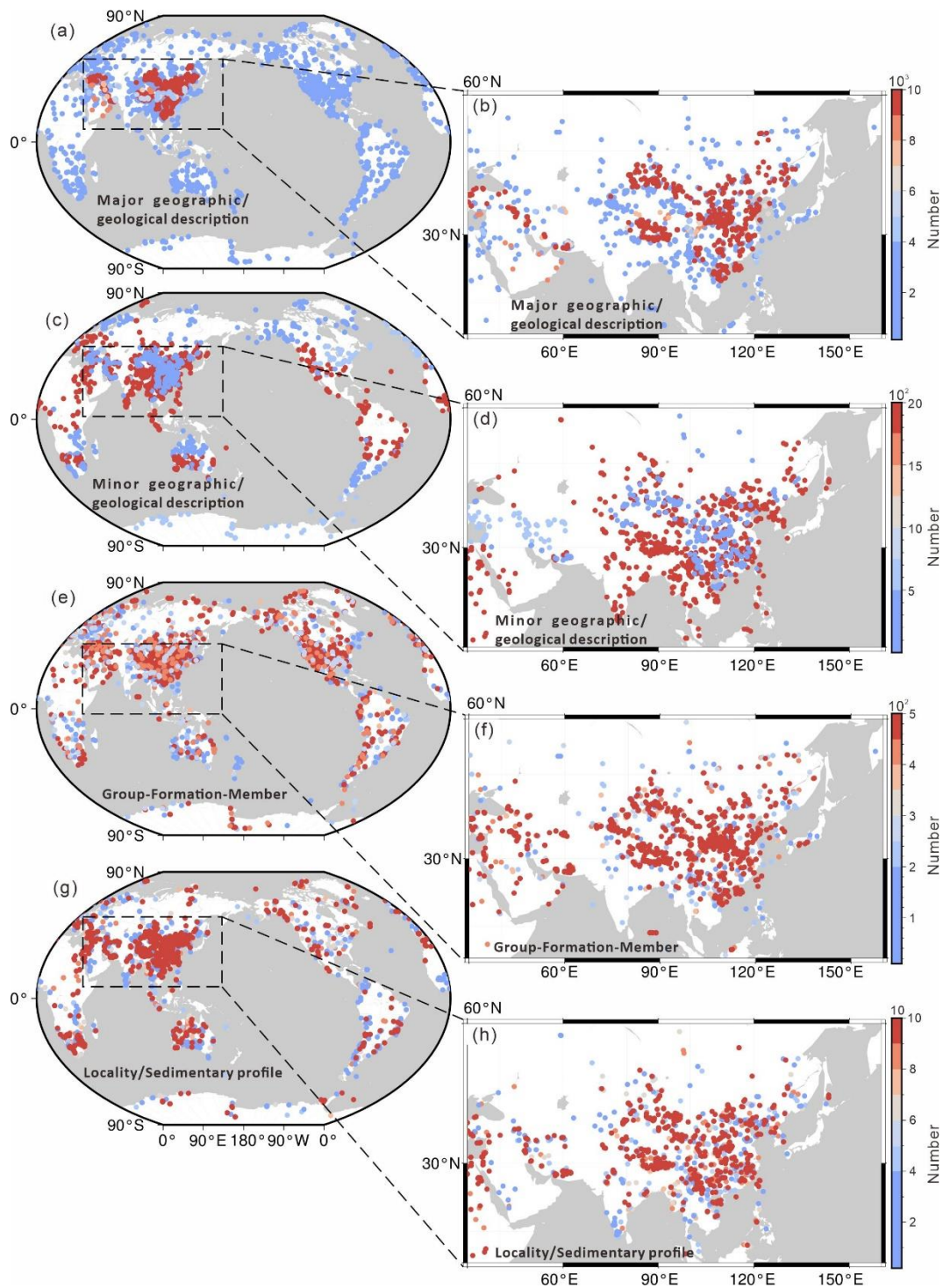
Figure 1: Workflow of constructing the OneDZ database (DataExpo was adopted from Lu et al., 2023, the DeepShovel tool was developed by Zhang et al., 2023, and the knowledge graph was based on Hu et al., 2024).



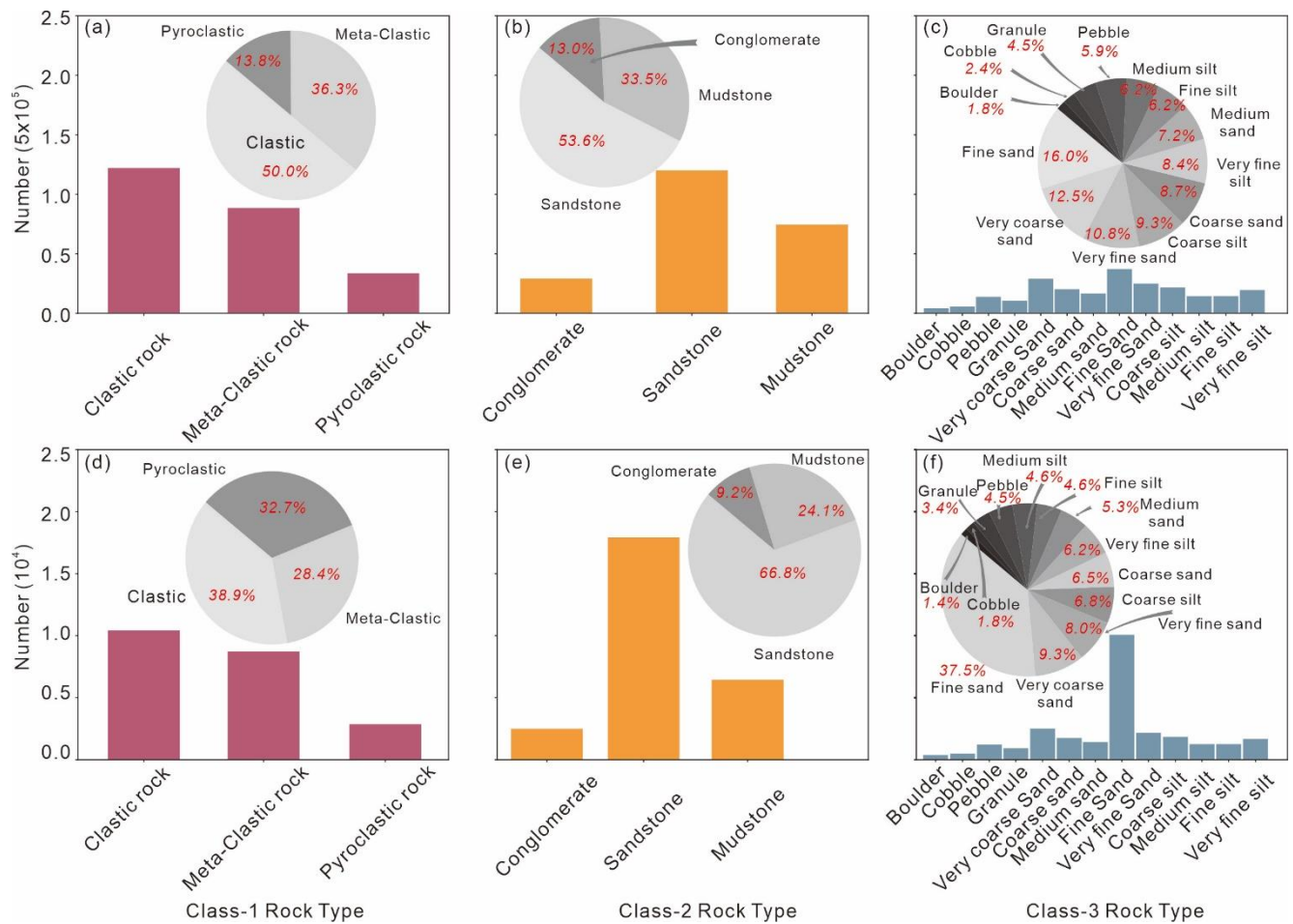


645

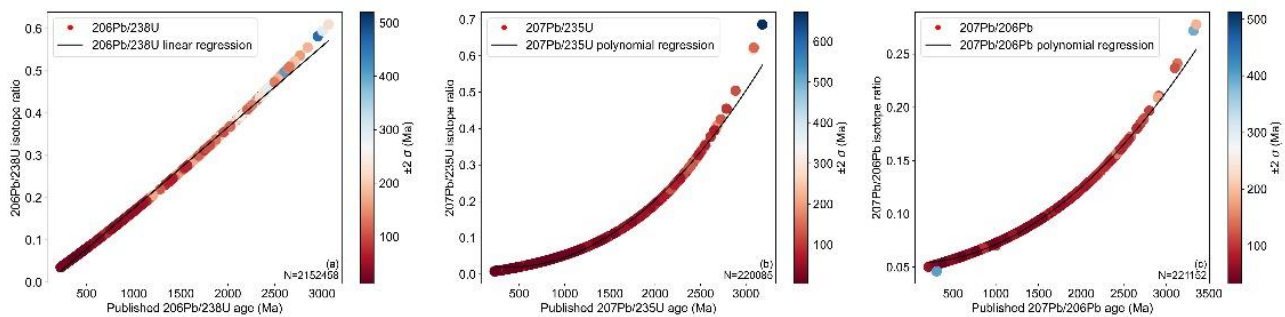
**Figure 3: Visualizations of the spatial, temporal and strata information in U-Pb dataset. (a)-(b) Major geographic/geological description; (c)-(d) Minor geographic/geological description; (e)-(f) Group-Formation-Member records; (g)-(h) Locality/Sedimentary profile.**



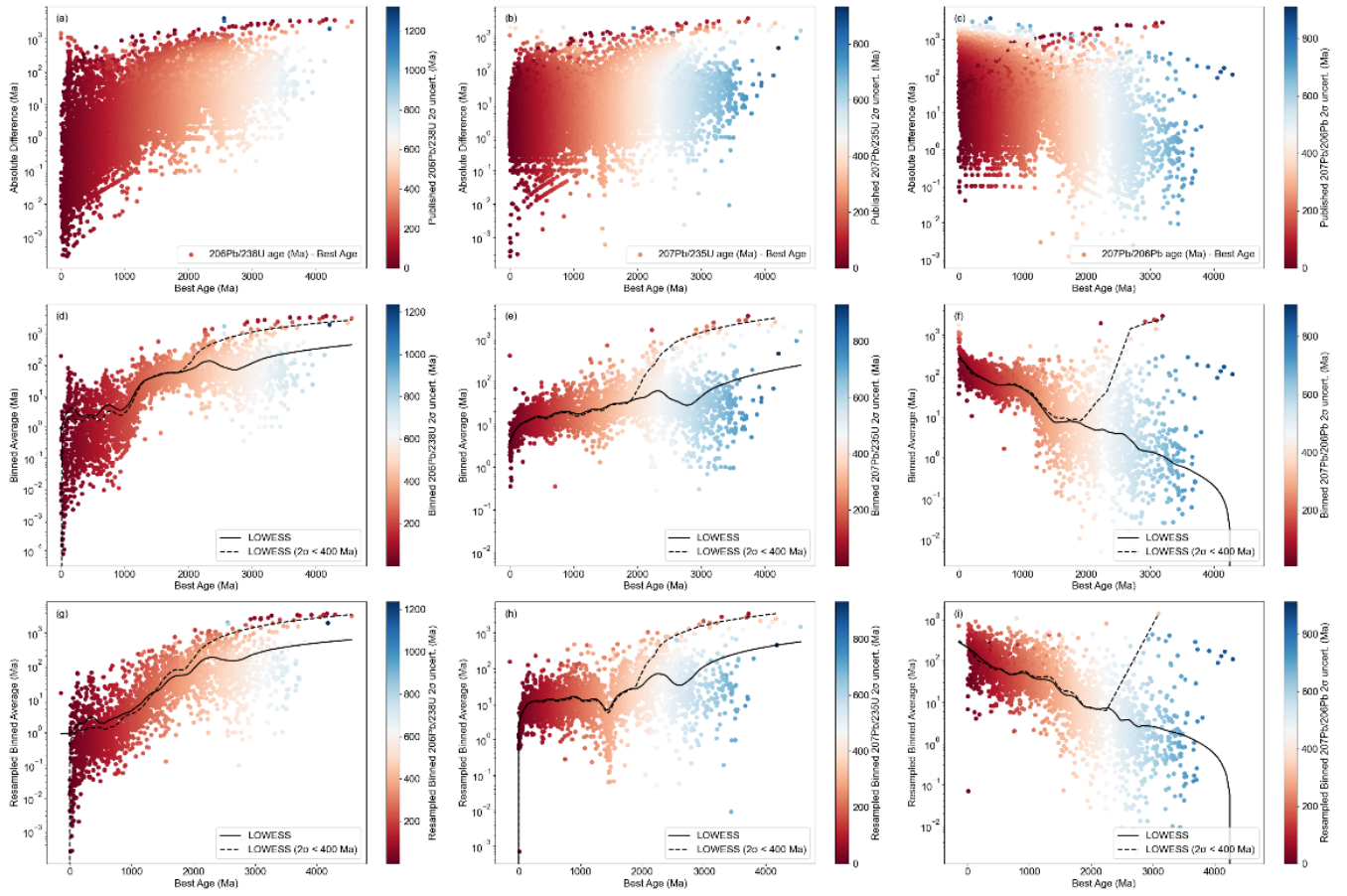
650 **Figure 4: Visualizations of the spatial, temporal and strata information in Lu-Hf dataset. (a)-(b) Major geographic/geological description; (c)-(d) Minor geographic/geological description; (e)-(f) Group-Formation-Member records; (g)-(h) Locality/Sedimentary profile.**



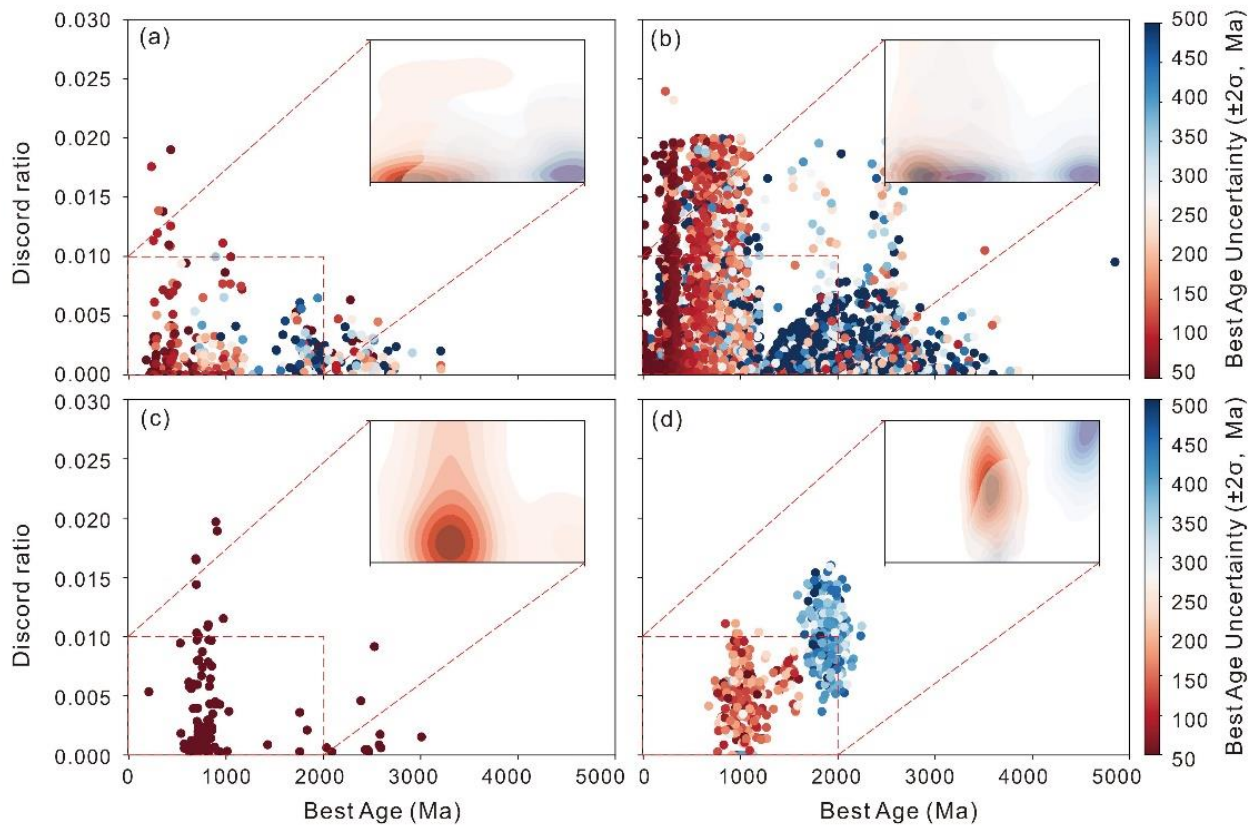
655 **Figure 5: Statistics of the rock types. (a) Class-1 type in U-Pb database; (b) Class-2 type in U-Pb database; (c) Class-3 type in U-Pb database; (d) Class-1 type in Lu-Hf database; (e) Class-2 type in Lu-Hf database; (f) Class-3 type in Lu-Hf database.**



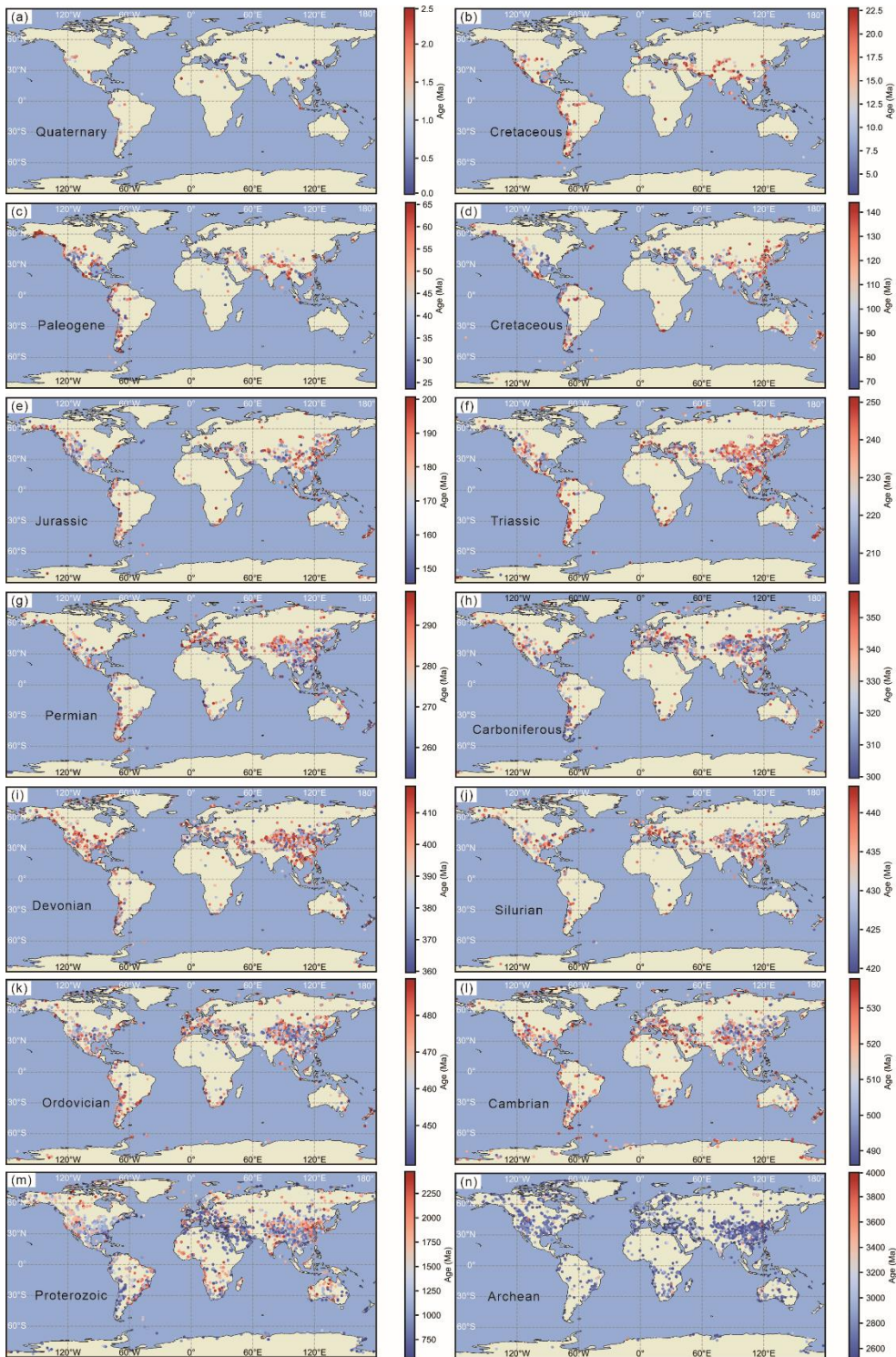
**Figure 6: Ages errors of different isotopic systems. (a)  $^{206}\text{Pb}/^{238}\text{U}$ ; (b)  $^{207}\text{Pb}/^{235}\text{U}$ ; (c)  $^{207}\text{Pb}/^{206}\text{Pb}$ .**



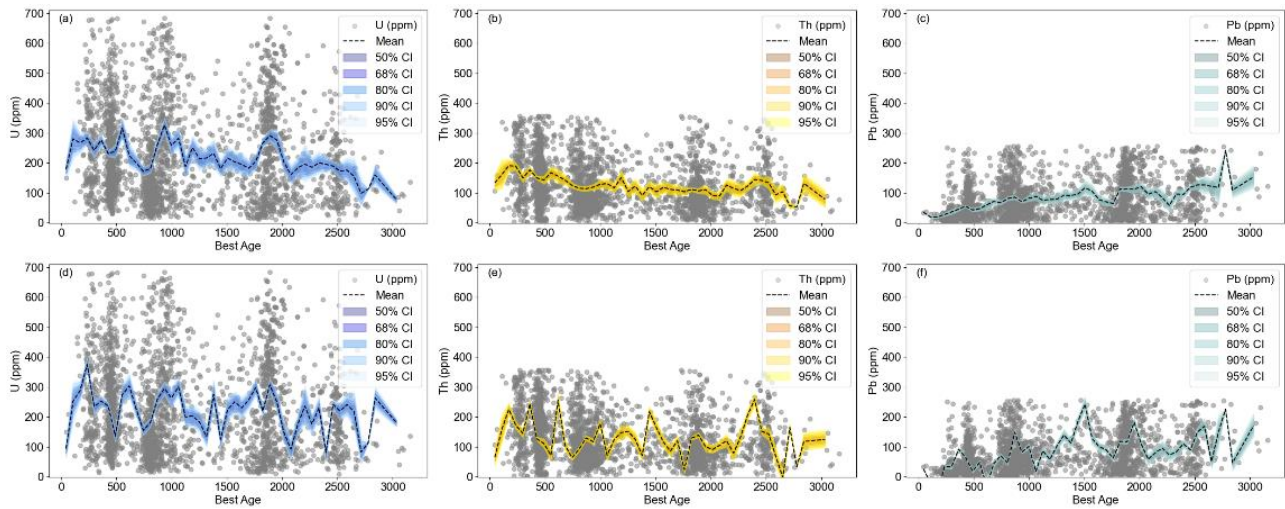
660 **Figure 7: Time-series of dating error via different isotopes. (a)-(c) Original data distribution; (d)-(f) Resampled by Monte-Carlo method; (g)-(i) Resampled by bootstrap method.**



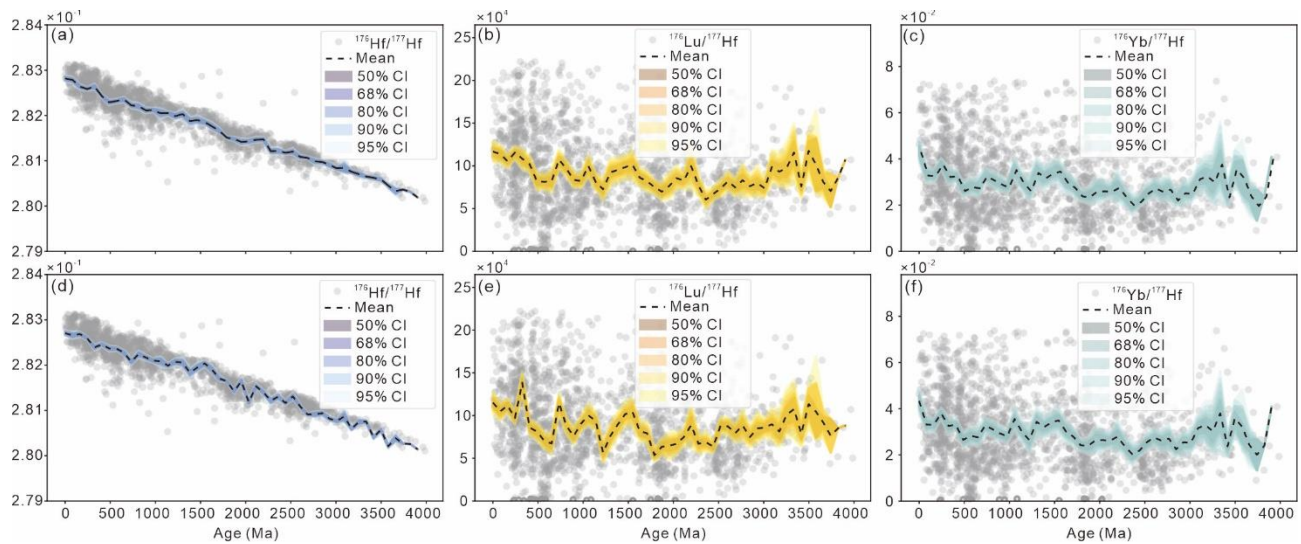
**Figure 8: Discordance ratio varying with time by different instruments. (a) SHRIMP; (b) LA-ICP-MS; (c) ID-TIMS; (d) SIMS.**



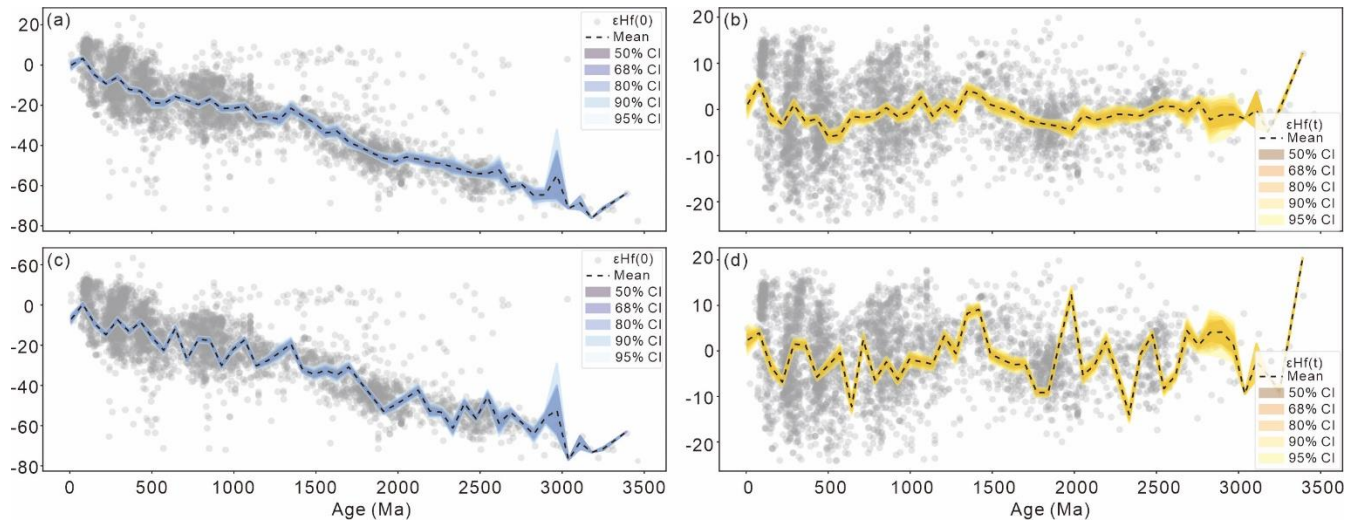
**Figure 9: Spatial-temporal distribution of the single U-Pb age record.**



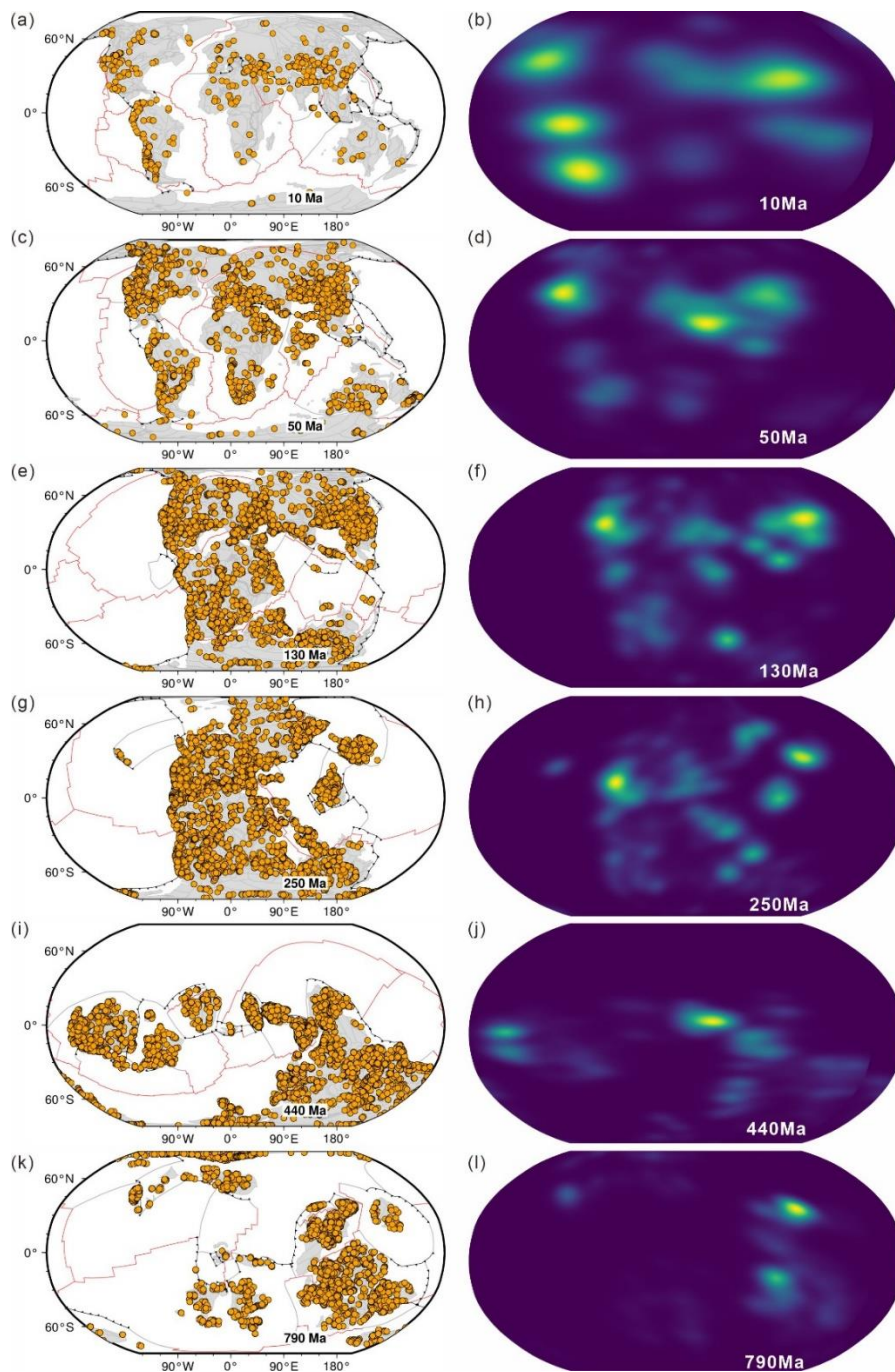
**Figure 10: Temporal variations of U, Th, Pb concentrations in U-Pb dataset. (a) U concentrations with bootstrap resampling; (b) Th concentrations with bootstrap resampling; (c) Pb concentrations with Monte-Carlo resampling; (d) U concentrations with Monte-Carlo resampling; (e) Th concentrations with Monte-Carlo resampling; (f) Pb concentrations with Monte-Carlo resampling..**



**Figure 11: Temporal variations of isotopic uncertainties in Lu-Hf dataset. (a)  $^{176}\text{Hf}/^{177}\text{Hf}$  with bootstrap resampling; (b)  $^{176}\text{Lu}/^{177}\text{Hf}$  with bootstrap resampling; (c)  $^{176}\text{Yb}/^{177}\text{Hf}$  with Monte-Carlo resampling; (d)  $^{176}\text{Hf}/^{177}\text{Hf}$  with Monte-Carlo resampling; (e)  $^{176}\text{Lu}/^{177}\text{Hf}$  with Monte-Carlo resampling; (f)  $^{176}\text{Yb}/^{177}\text{Hf}$  with Monte-Carlo resampling.**



675 **Figure 12: Temporal variations of  $\epsilon_{\text{Hf}}$  uncertainties in Lu-Hf dataset. (a)  $\epsilon_{\text{Hf}}(0)$  with bootstrap resampling; (b)  $\epsilon_{\text{Hf}}(t)$  with bootstrap resampling; (c)  $\epsilon_{\text{Hf}}(0)$  with Monte-Carlo resampling; (d)  $\epsilon_{\text{Hf}}(t)$  with Monte-Carlo resampling.**



**Figure 13: Paleo-distributions and spatial kernel density estimate of U-Pb records (the tectonic model was from Merdith et al., 2021 and the temporal resolution is  $1^\circ \times 1^\circ$ ). (a)-(b) Paleo-distribution and density of 10Ma; (c)-(d) Paleo-distribution and density of 50Ma; (e)-(f) Paleo-distribution and density of 130Ma; (g)-(h) Paleo-distribution and density of 250Ma; (i)-(j) Paleo-distribution and density of 440Ma; (k)-(l) Paleo-distribution and density of 790Ma.**

680

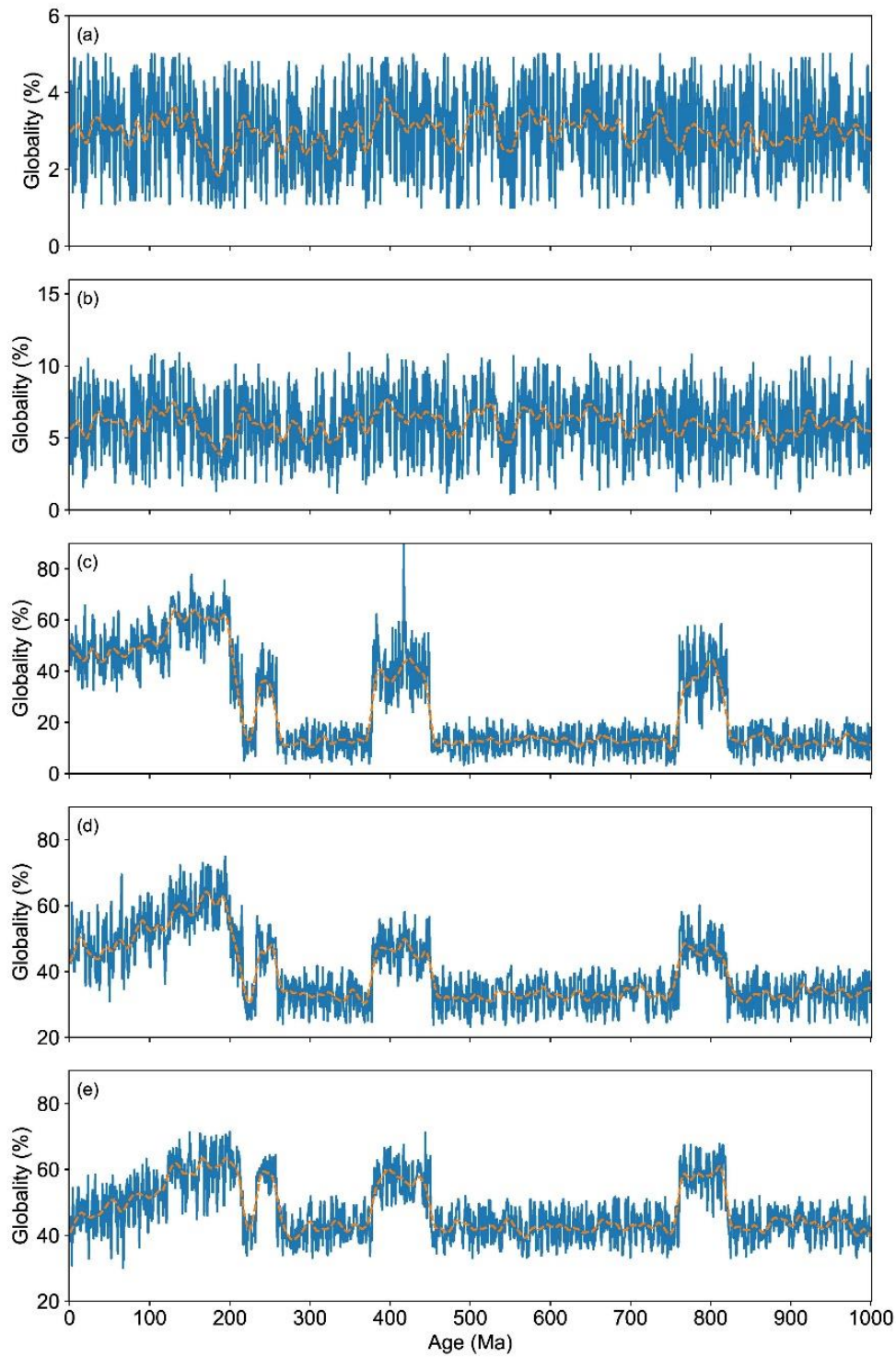
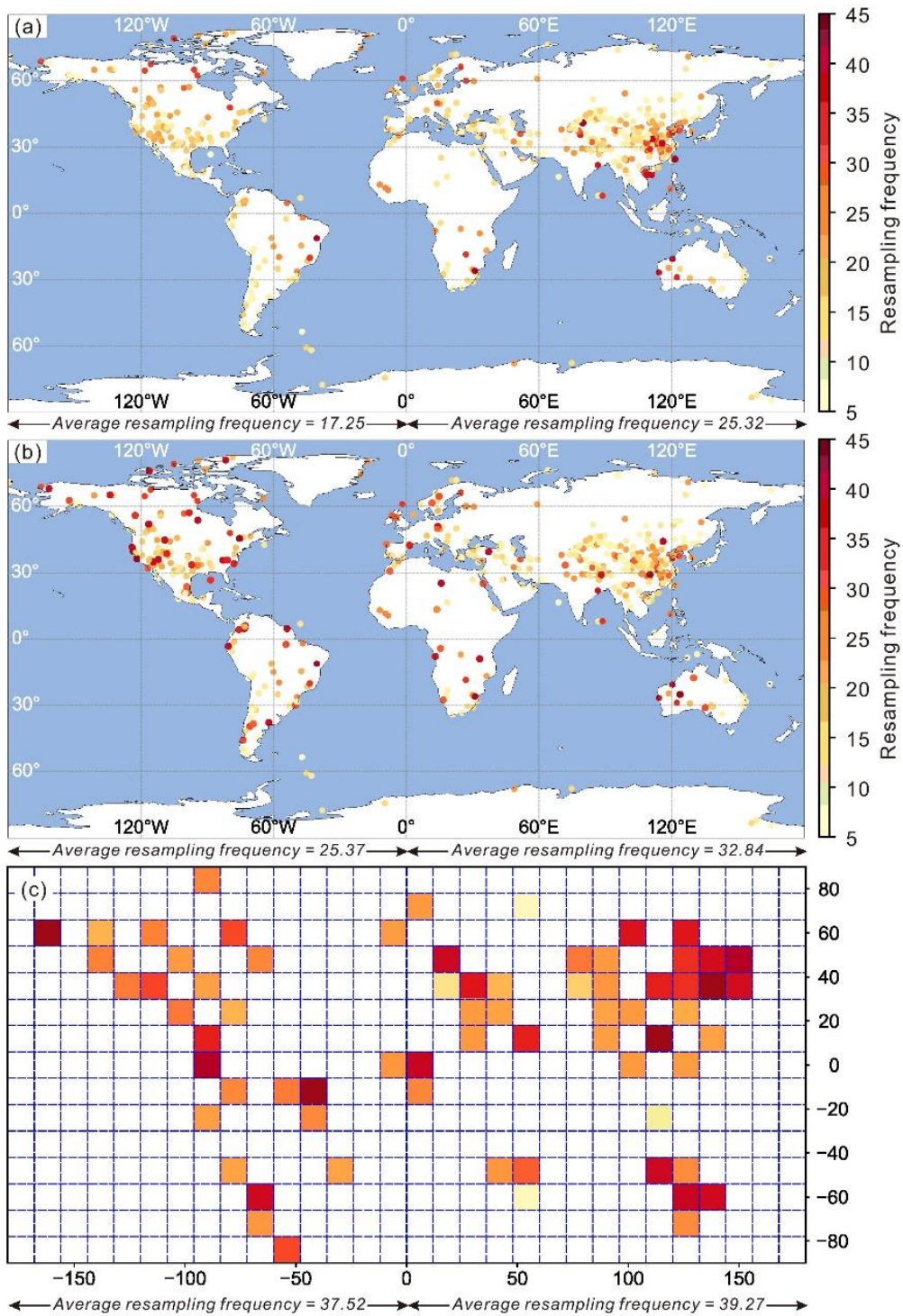


Figure 14: Global evaluation of U-Pb data with different grid sizes. (a) 2°; (b) 4°; (c) 6°; (d) 8°; (e) 10°.



685 **Figure 15: The resampling frequency of different methods. (a) Monte-Carlo method; (b) SMOTE-Monte-Carlo method; (c) 12°×12° grid-SMOTE-Monte-Carlo method.**

**Table 1: Construction methods comparison of three typical zircon databases**

Dataset	Methods	Information types	Data cleaning	Adapting AI tool
OneDZ	Crowdfunding	Data origin, Spatial information, Stratigraphic information, Isotopic information	Python and MySQL scripts and artificial checking	Yes
Wu et al., 2024	Directly collecting	Data origin, Spatial information, Isotopic information	Artificial checking	Not mentioned
Puetz et al., 2024	Directly collecting	Data origin, Spatial information, Stratigraphic information, Isotopic information	Artificial checking	Not mentioned

**Table 2: Data comparison of three typical zircon databases**

Dataset	Field number of U-Pb	Field number of Lu-Hf	Reference (10 <sup>4</sup> )	Sample (10 <sup>4</sup> )	Region	Geological unit	GPS (10 <sup>6</sup> )
OneDZ	<b>71</b>	<b>86</b>	<b>5.4</b>	<b>31</b>	<b>215</b>	<b>1348</b>	<b>1.8</b>
Wu et al., 2024	24	/	3.4	2.8	208	1347	0.5
Puetz et al., 2024	34	26	4.2	20	215	1305	0.6

**Note: the bolded number represents the largest number in different items. Only statistic the detrital zircon from Wu et al., 2024 and Puetz et al., 2024.**

690

**Table 3: Data specifications of the reference information (the proportion was calculated by number of valid items divided the number of total items)**

Field Name	Description	Proportion in U-Pb data		Proportion in Lu-Hf data	
		Total	After expert checking	Total	After expert checking
Lead_Author	Information of the authors	96.03%	41.56%	100.00%	100.00%
Year	The year the paper published	96.02%	41.56%	100.00%	100.00%
Journal	The journal the paper the published	95.56%	41.56%	100.00%	100.00%
Vol	Volume	92.90%	41.56%	100.00%	100.00%
Pages	Pages	91.91%	41.39%	100.00%	100.00%
Title	Paper's title	96.04%	41.56%	99.98%	99.98%
Web_Link	The link of the published paper	53.16%	20.61%	93.95%	93.95%

**Table 4: Data specifications of the sample, spatial and strata information (the proportion was calculated by number of valid items divided the number of total items)**

Field Name	Description	Proportion in U-Pb data		Proportion in Lu-Hf data	
		Total	After expert checking	Total	After expert checking
Published_Sample_ID	Sample ID from original publication	79.45%	41.56%	99.41%	99.41%
Country_State	Country or state/province	74.91%	41.56%	93.12%	93.12%
Region	Geographic region	79.75%	41.56%	98.53%	98.53%
Continent	Continent	91.87%	41.56%	91.91%	91.91%
Major_Geographic_Geologic_Unit	Major geologic or geographic unit	74.77%	41.56%	98.88%	98.88%
Minor_Geologic_Geographic_Unit	Minor geologic or geographic unit	67.96%	41.53%	95.45%	95.45%
Group	Stratigraphic group	11.64%	0.19%	4.94%	4.94%
Formation	Stratigraphic formation	23.57%	0.44%	7.65%	7.65%
Member	Stratigraphic member	3.36%	0.25%	2.22%	2.22%
Locality	Sampling locality	81.81%	40.99%	81.26%	81.26%
Profile	Stratigraphic section or drill profile	4.05%	0.20%	11.54%	11.54%
Latitude	Latitude (decimal degrees)	95.08%	41.56%	97.08%	97.08%
Longitude	Longitude (decimal degrees)	95.08%	41.56%	96.06%	96.06%
Depos_Age_Period	Depositional period	12.76%	0.74%	12.28%	12.28%
Depos_Age_Epoch	Depositional epoch	9.01%	0.56%	8.09%	8.09%
Depos_Age_Stage	Depositional stage	1.20%	0.23%	3.92%	3.92%
Max_Depos_Age_Ma	Maximum depositional age (Ma)	48.85%	40.82%	8.85%	8.85%
Est_Depos_Age_Ma	Estimated depositional age (Ma)	49.33%	40.82%	46.17%	46.17%
Min_Depos_Age_Ma	Minimum depositional age (Ma)	30.63 %	23.5%	6.99%	6.99%

**Table 5: Data specifications of the U-Pb isotopic system (the proportion was calculated by number of valid items divided the number of total items)**

Field Name	Description	Proportion in total dataset	Proportion after expert checking
Spectrometer	Mass spectrometer	41.56%	88.28%
Spectrometer_Location	Laboratory location	24.46%	55.74%
Institution	Analytical institution	24.46%	58.03%
Grain	Zircon grain identifier	41.56%	93.66%
Spot_Location	Ablation spot position	40.82%	75.95%
Spot_diam	Spot diameter ( $\mu\text{m}$ )	24.46%	31.34%
Pb206U238_iso	$^{206}\text{Pb}/^{238}\text{U}$ isotopic ratio	24.24%	70.91%
Pb206U238_iso_one_sigma	$^{206}\text{Pb}/^{238}\text{U}$ ratio $1\sigma$ uncertainty	24.24%	66.96%
Pb207U235_iso	$^{207}\text{Pb}/^{235}\text{U}$ isotopic ratio	24.24%	45.95%
Pb207U235_iso_one_sigma	$^{207}\text{Pb}/^{235}\text{U}$ ratio $1\sigma$ uncertainty	24.24%	47.22%
Pb207Pb206_iso	$^{207}\text{Pb}/^{206}\text{Pb}$ isotopic ratio	24.21%	46.61%
Pb207Pb206_iso_one_sigma	$^{207}\text{Pb}/^{206}\text{Pb}$ ratio $1\sigma$ uncertainty	24.21%	43.79%
Pb208Th232_iso	$^{208}\text{Pb}/^{232}\text{Th}$ isotopic ratio	0.13%	3.77%
Pb208Th232_iso_one_sigma	$^{208}\text{Pb}/^{232}\text{Th}$ ratio $1\sigma$ uncertainty	0.13%	3.76%
Pb206U238_age	$^{206}\text{Pb}/^{238}\text{U}$ age (Ma)	22.25%	97.20%
Pb206U238_age_one_sigma	$^{206}\text{Pb}/^{238}\text{U}$ age $1\sigma$ uncertainty	0.25%	66.96%
Pb206U238_age_two_sigma	$^{206}\text{Pb}/^{238}\text{U}$ age $2\sigma$ uncertainty	40.37%	80.85%
Pb207U235_age	$^{207}\text{Pb}/^{235}\text{U}$ age (Ma)	20.01%	90.15%
Pb207U235_age_one_sigma	$^{207}\text{Pb}/^{235}\text{U}$ age $1\sigma$ uncertainty	0.25%	51.15%
Pb207U235_age_two_sigma	$^{207}\text{Pb}/^{235}\text{U}$ age $2\sigma$ uncertainty	36.84%	73.79%
Pb207Pb206_age	$^{207}\text{Pb}/^{206}\text{Pb}$ age (Ma)	38.61%	92.84%
Pb207Pb206_age_one_sigma	$^{207}\text{Pb}/^{206}\text{Pb}$ age $1\sigma$ uncertainty	0.25%	52.11%
Pb207Pb206_age_two_sigma	$^{207}\text{Pb}/^{206}\text{Pb}$ age $2\sigma$ uncertainty	38.36%	76.16%
Best_Age	Preferred best age (Ma)	18.03%	69.98%
Best_age_one_sigma	Best age $1\sigma$ uncertainty	0.72%	27.25%
Best_age_two_sigma	Best age $2\sigma$ uncertainty	17.80%	33.10%
Discord	Discordance (%)	0.74%	25.53%
U_ppm	Uranium concentration (ppm)	0.31%	13.10%
Th_ppm	Thorium concentration (ppm)	0.23%	9.88%
Pb_ppm	Lead concentration (ppm)	0.11%	5.40%
ThU_ratio	Th/U atomic ratio	0.53%	23.49%

**Table 6: Data specifications of the Lu-Hf isotopic system (the proportion was calculated by number of valid items divided the number of total items)**

Field Name	Description	Proportion in total dataset	Proportion after expert checking
Spectrometer	Mass spectrometer	11.46%	11.46%
Spectrometer_Location	Laboratory location	10.58%	10.58%
Institution	Analytical institution	10.37%	10.37%
Grain	Zircon grain identifier	19.17%	19.17%
Spot_Location	Ablation spot position	4.55%	4.55%
Spot_diam	Spot diameter ( $\mu\text{m}$ )	4.52%	4.52%
Upb_Age	U-Pb age (Ma)	95.02%	95.02%
Upb_Age_two_sigma	U-Pb age $2\sigma$ uncertainty	10.74%	10.74%
U_Pb_Age_uncertainty_2sigma	U-Pb age absolute $2\sigma$ uncertainty	13.02%	13.02%
176Hf177Hf_iso	$^{176}\text{Hf}/^{177}\text{Hf}$ isotopic ratio	95.06%	95.06%
176Hf177Hf_iso_2sigma	$^{176}\text{Hf}/^{177}\text{Hf}$ ratio $2\sigma$ uncertainty	87.09%	87.09%
176Lu177Hf_iso	$^{176}\text{Lu}/^{177}\text{Hf}$ isotopic ratio	95.48%	95.48%
176Lu177Hf_iso_2sigma	$^{176}\text{Lu}/^{177}\text{Hf}$ ratio $2\sigma$ uncertainty	56.65%	56.65%
176Yb177Hf_iso	$^{176}\text{Yb}/^{177}\text{Hf}$ isotopic ratio	72.64%	72.64%
176Yb177Hf_iso_2sigma	$^{176}\text{Yb}/^{177}\text{Hf}$ ratio $2\sigma$ uncertainty	46.91%	46.91%
epsilon_Hf_0	$\epsilon\text{Hf}(0)$ value	13.01%	13.01%
epsilon_Hf_0_1sigma	$\epsilon\text{Hf}(0)$ $1\sigma$ uncertainty	2.89%	2.89%
epsilon_Hf_0_2sigma	$\epsilon\text{Hf}(0)$ $2\sigma$ uncertainty	2.58%	2.58%
epsilon_Hf_t	$\epsilon\text{Hf}(t)$ value	86.55%	86.55%
epsilon_Hf_t_1sigma	$\epsilon\text{Hf}(t)$ $1\sigma$ uncertainty	8.24%	8.24%
epsilon_Hf_t_2sigma	$\epsilon\text{Hf}(t)$ $2\sigma$ uncertainty	75.43%	75.43%
TDM1_Ma	Hf single-stage model age (Ma)	77.76%	77.76%

TDM1_Ma_2sigma	TDM1 2 $\sigma$ uncertainty	1.86%	1.86%
TDM2_Ma	Hf two-stage model age (Ma)	92.24%	92.24%
TDM2_Ma_2sigma	TDM2 2 $\sigma$ uncertainty	1.47%	1.47%

---

700