

5



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

Keran Li¹, Xiumian Hu^{1,} *, Rong Chai², Jianghai Yang³, Weiwei Xue⁴, Yingdi Pan¹, Taiyang Li⁵, Can Fang⁶, Anlin Ma¹, Hu Huang^{7,8}, Qianqian Guo⁹, Wentao Yang¹⁰, Lisha Hu¹¹, Liang Qi^{7,8}, Guohui Chen¹², Gaoyuan Sun¹³, Shijie Zhang¹⁴, Tao Deng¹, Kuizhou Li^{7, 15}, Jiaopeng Sun¹⁶, Biao Gao¹⁷

¹State Key Laboratory of Mineral Deposit Research, School of Earth Sciences and Engineering, Nanjing University, Nanjing 210023, China
 ²Chinese Academy of Geological Sciences, Beijing 100037, China
 ³School of Earth Sciences, China University of Geosciences (Wuhan), Wuhan 430074, China
 ⁴State Key Laboratory of Isotope Geochemistry, Guangzhou Institute of Geochemistry, Chinese Academy of Sciences (CAS),

⁴State Key Laboratory of Isotope Geochemistry, Guangzhou Institute of Geochemistry, Chinese Academy of Sciences (CAS), Guangzhou, China
⁵College of Computer Science and Cyber Security, Chengdu University of Technology, Chengdu 610059, China

⁶Hangzhou Research Institute, Huawei Technologies, Hangzhou 310056, China

⁷State Key Laboratory of Oil and Gas Reservoir Geology and Exploitation, Institute of Sedimentary Geology, Chengdu
University of Technology, Chengdu 610059, China

- ⁸Key Laboratory of Deep-time Geography and Environment Reconstruction and Applications of Ministry of Natural Resources, Chengdu University of Technology, Chengdu 610059, China
 ⁹College of Earth and Planetary Sciences, University of Chinese Academy of Sciences, Beijing, 100049, China
 ¹⁰School of Resources and Environment, Henan Polytechnic University, Jiaozuo 454000, China
- ¹¹College of Marine Geosciences, Ocean University of China, Qingdao 266100, China
 ¹²School of Earth Sciences and Engineering, Hohai University, Nanjing 210098, China
 ¹³College of Oceanography, Hohai University, Nanjing 210024, China
 ¹⁴College of Tourism, Henan Normal University, Xinxiang, China
 ¹⁵College of Earth and Planetary Sciences, Chengdu University of Technology, Chengdu 610059, China
- ¹⁶State Key Laboratory of Continental Dynamics, Department of Geology, Northwest University, Xi'an 710069, China
 ¹⁷State 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 with the continuous

- 30 improvement of testing 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 users. Comparing with
- 35 current zircon database, OneDZ complies 1,925,687 gains of detrital zircon U-Pb and 275,971 gains 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 time. It serves as a springboard for research, offering new insights in understanding Earth's past, present, and future. The database (Li and Hu, 2025) is freely available via Zenodo at https://doi.org/10.5281/zenodo. 15522949. All code snippets in this research are accessible via https://github.com/KeranLi/Global-Detrital-Zircon. The OneDZ web platform is accessible via https://dedc.geoscience.cn/onedz/.

1 Introduction

- 45 The advent of high-precision and high-accuracy U-Pb geochronology has revolutionized the understanding of Earth's history. Since the development of isotope dilution thermal ionization mass spectrometry (ID-TIMS) (Krogh, 1973) and subsequent advancements in mass spectrometry techniques (Jarvis and Kym, 1988; MacRae and Neil, 1995; Belu et al., 2003; Muzikar et al., 2003; Yergey et al., 2013), the ability to date zircon crystals with exceptional precision has significantly improved. Zircon, a robust and ubiquitous mineral found throughout the continental crust, serves as a reliable recorder of geological events due
- 50 to its high closure temperature and resistance to weathering and metamorphism (Pupin, 1980). The magmatic or metamorphic origin zircons are often broken, transported, and stored in sediments or sedimentary rocks, known as detrital zircons. The analysis for detrital zircons includes U-Pb and Lu-Hf isotopic systems. Chemical formula of detrital zircon can be represented as [ZrSiO₄]. The ionic radius of [Zr⁴⁺] is 0.87 Å, which can be easily replaced by [U⁴⁺] and [Th⁴⁺] because of
- similar ionic radius of 1.05 Å and 1.10 Å (Jaffey et al., 1971). Two isotopes of $[U^{4+}]$ (²³⁸U and ²³⁵U) generate ²⁰⁶Pb and ²⁰⁷Pb 55 isotopes following the decay processes: ²³⁸U \rightarrow ²⁰⁶Pb+8 α +6 β ⁻ (half-life: 4468 million years, Jaffey et al., 1971), ²³⁵U \rightarrow ²⁰⁷Pb+7 α +4 β ⁻ (half-life: 703.8 million years, Jaffey et al., 1971) and ²³²Th \rightarrow ²⁰⁸Pb+6 α +4 β ⁻ (half-life: 1400 million years, Jaffey et al., 1971). Based on triple decay processes, the detrital zircon ages can be obtained via consisted ²⁰⁶Pb/²³⁸U, ²⁰⁷Pb/²³⁵U and ²⁰⁸Pb/²³²Th decayed ages.

In addition to U-Pb geochronology, the Lu-Hf isotopic system has become an indispensable tool for understanding crustal

- 60 evolution and mantle differentiation (Patchett and Tatsumoto, 1983). The [Lu³⁺] is the heaviest rare earth element (REE) and are easily enriched in detrital zircon. The ¹⁷⁶Lu decays to ¹⁷⁶Hf via ¹⁷⁶Lu \rightarrow ¹⁷⁶Hf+ β ⁻ (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 (Cherniak et al., 1997). The Lu-Hf isotopic data are noted by ϵ units by ϵ Hf(0)=10000×[(¹⁷⁶Hf/¹⁷⁷Hf)_{sample}/(¹⁷⁶Hf/¹⁷⁷Hf)_{CHUR,0}-1] and ϵ Hf(t)=10000×{[(¹⁷⁶Hf/¹⁷⁷Hf)_{sample}-(¹⁷⁶Lu/¹⁷⁷Hf)_{sample}×(e^{λ t}-1)]/[(¹⁷⁶Hf/¹⁷⁷Hf)_{CHUR},0-(¹⁷⁶Lu/¹⁷⁷Hf)_{CHUR}×(e^{λ t}-1)]-1}. t is the
- 65 crystallization age. ${}^{176}\text{Hf}/{}^{177}\text{Hf}$ and ${}^{176}\text{Lu}/{}^{177}\text{Hf}$ can be measured from detrital zircons. λ is the decay constant and equals to 1.867×10^{-5} million years (Söderlund et al., 2004). CHUR denotes the isotopic results of the chondritic uniform reservoir. In the past two decades, it is estimated millions or even more U-Pb geochronological data of detrital zircons internationally have been reported. As the amount of data increases, it's possible for using 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
- 70 of continental arcs (McKenzie et al., 2016; Cao et al., 2017), tectonic history (Cawood et al., 2012; Barham et al., 2022; Zhang



Science Science

et al., 2023; Malone et al., 2024; Odlum et al., 2024), crust evolution (Cheng, 2017; Barham et al., 2019; Cawood, 2020), paleo-geographic (Xue et al., 2022; Jian et al., 2022) and provenance analysis (Wang et al., 2024). Along with data-driven analysis, several analysis 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,

75 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 effectively cover detrital zircon data in sedimentary rocks.

In addition, the current database construction mainly focuses on reporting data, lacking discussion on the problems that exist in the database construction process. Several previous reported databases tend to store data by Microsoft Excel software (Voice

- 80 et al., 2011; Puetz, 2019; Martin et al., 2022; Puetz, 2024; Wu et al., 2024). However, the Excel software can only store limited raw data (n=1048576). Some researches split the data sheets into several Excel files like Wu et al. (2022). The split data sheets hind fast data searching and the huge amount of data results in excessively low access efficiency in Excel software. Therefore, the construction of giant earth science databases represented by detrital zircon data urgently needs to shift towards the use of more professional database software. Such curated compilations, like Earth Chemistry (EarthChem) and Geochemistry of
- 85 Rocks of the Oceans and Continents (GEOROC), have been developed. These platforms all utilize database software and frontend and back-end to retrieve and download data on web pages. Convenient operation on the web page will guarantee continuously data updating, providing powerful and sustainable data sources for future research. In order to address the challenges in current geological research based on detrital zircon big data, we have established a detrital

zircon database that covers both English and Chinese literature worldwide. Presented here is the OneDZ database, an extensive

- 90 compilation of zircon U-Pb geochronological and Lu-Hf isotopic data, encompassing over 1925687 U-Pb and 275971 Lu-Hf records from approximately 275971 publications. This database 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 dating instruments, host rock lithologies, stratigraphic information, and other original records. OneDZ records the lithology, stratigraphic, spatial, and testing information of detrital zircons as much as possible. In the construction of OneDZ database,
- 95 the Excel software was abandoned, and instead, professional database management software such as MySQL was adapted. In addition, several AI and code snippets were introduced and displayed the improvement of artificial intelligence tools in database construction. In order to make it more convenient for using, OneDZ has completed the deployment on the web side. Global and regional interdisciplinary research can be conducted simultaneously on OneDZ. At the same time, the enormous data also makes OneDZ a natural laboratory for discussing data analysis methods in Earth science. OneDZ provides a
- 100 foundation for research in multiple aspects, including data provision, database construction, and discussion and analysis of data analysis methods in earth science.



2 Database construction

One of the most different features in OneDZ is the systematic construction workflow (Fig. 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. 105 With knowledge graph, the words to describe the sample location, sedimentary or stratigraphic descriptions and the isotopic results are most tight information associated to detrital zircon researches. Although almost no previous research summarized the difficulties in collecting data sources, 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 research, AI-assisted tools, including DataExpo and GPT Agent (Supplement 1), were employed to check specific online resources like

110 Pangea (https://pangaea.de/), Google Scholar, and CNKI to search potential papers containing data. 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. Those validated data were incorporated into the OneDZ database. Details of construction differences are shown in Table 1.

115 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

- 120 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
- 125 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).

2.2 Facility from AI tools

One of the fundamental challenges in constructing a database lies in data collection. Although the crowdfunded approach ensures the geologists participating in database development are experts in the research area, their expertise does not guarantee

130 familiarity with every publication. To ensure no potential metadata is missing, this study introduced a data cruise system integrated with deep learning technology. Named as DataExpo (Lu et al., 2023), the cruise system employs deep learning for metadata extraction (Figure S1-S2 in the Supplement), performing automatic semantic tagging, classification, and structured



140

information extraction from web pages. DataExpo automatically crawls web pages related to detrital zircon. Using a multidimensional web page ranking strategy, retrieval results for different queries are sorted. Finally, based on natural language

135 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, a GPT 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.

2.3 Automatic data process

- 145 Existing databases on zircon and other earth sciences typically emphasize data quality and related aspects. However, in the current era of data growth, building specialized databases in earth sciences will become increasingly common, but few studies have focused on the critical aspect of data cleaning. 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
- 150 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 deduplication. Following the standard data cleaning process, Python and MySQL scripts were designed for detecting missing key items, checking for conflicting content, detecting format anomalies, and Supplementing duplicate data entries (see Supplement 2 for details). Python scripts are more flexible while MySQL scripts have higher running efficiency (Figure S4 in 155 the Supplement).
- 155 the Supplement).

3 Database

In our OneDZ database, 1925687 detrital zircon U-Pb age data points and 275971 Lu-Hf isotope analyses were compiled. From multiple dimensions such as region, literature, and samples, OneDZ is currently the most comprehensive database for global detrital zircon data records (Table 2). The U-Pb geochronological data cover potential research regions (Figure. 2a).

160 The Lu-Hf data are primarily distributed across China, South Africa, India, and Australia (Figure. 2b). Periodic statistics indicate that ancient zircons (over 1000 Ma) predominantly contribute to this database in both U-Pb and Lu-Hf data (Figure. 2c-f). The specific data situation of different fields in OneDZ can be seen in Table 3-6.



3.1 Reference information

- 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
- 170 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). In Lu-Hf data, the same journals also mostly contribute to OneDZ (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
- 175 diversity (Figure S6a-b in the Supplement).

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

- 180 research. Due to the fact that geographic information Supplementation is not a must item in databases, high-precision latitude and longitude coordinates are still the preferred choice for spatial information. 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
- 185 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 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).
- 190 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 (Fig. 3) and Lu-Hf (Fig. 4) datasets continues to exhibit significant spatial skew. A majority of the records are concentrated in East Asia, with a particular focus on China. This concentration is not solely due to oversampling but also reflects a propensity towards more extensive research coverage in these regions.

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 data are accompanied by detailed instrumental information, including the types of mass spectrometer utilized (e.g., LA-ICP-MS, SHRIMP, and ID-TIMS), the analysis institution's spot types (rim or core), and the spot diameter measurements. For the chronological data, the isotopic ratios ²⁰⁶Pb/²³⁸U, ²⁰⁷Pb/²³⁵U, ²⁰⁷Pb/²⁰⁶Pb, and ²⁰⁸Pb/²³²Th were recorded

with corresponding 1σ uncertainties. A limited number of papers have reported uncertainties at the 2σ level. OneDZ complied all uncertainties following the original data. In cases where the preferred age is not explicitly stated, the most reliable age estimate, considering 1σ and 2σ uncertainties, is selected with predefined thresholds of 1200 Ma and 1600 Ma, following the guidelines established by Gehrels et al., 2008.

Furthermore, the database also archives the discord ratio, concentrations of U, Th, and Pb, as well as the U/Th and 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 Yb/ 177 Hf, 176 Lu/ 177 Hf, and 176 Hf/ 177 Hf isotopic ratios, each accompanied by their corresponding 2σ uncertainties, which reflect the precision of the

215 measurements.

210

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

220 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,



the OneDZ database collected petrological contexts from original articles. OneDZ categorized rock types into a hierarchical
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. 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 5b). As for the Class-3 rock types, the proportion of rocks
with different particle sizes is relatively close, with fine sand and very coarse sand being the most common types (Figure 5c).
In the Lu-Hf datasets, relatively little rock records were provided in the original article. The clastic rock majorly contributed
to the datasets (38.9%, 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). As for the Class-3 rock types, the proportion of rocks with different particle sizes is relatively close, with fine sand being the most common types (37.5%, Figure 5f).

4.2 Data uncertainty

The data uncertainty in 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-

- 240 lives among different isotopic systems. Dating uncertainties and potential biases are more concerned about data processing. Figure 6 illustrates the relationships between isotopic ratios, calculated ages, and their corresponding 2σ uncertainties. The ²⁰⁶Pb/²³⁸U 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 2σ uncertainty of ages escalates to around 300 Ma. This trend suggests that approximately 67% of samples may be associated
- 245 with a temporal uncertainty of approximately 600 Ma. In contrast, the ²⁰⁷Pb/²³⁵U and ²⁰⁷Pb/²⁰⁶Pb isotopic systems are characterized by second-order polynomial regressions (Figure 6b-c). The complex regression models suggest a 2σ uncertainty of 500 Ma emerging at around 3000 Ma in ²⁰⁷Pb/²³⁵U and ²⁰⁷Pb/²⁰⁶Pb isotopic systems. The age uncertainty becomes significantly pronounced when analyzing samples over 3000 Ma. The relatively delayed error accumulation effect renders the ²⁰⁷Pb/²³⁵U and ²⁰⁷Pb/²⁰⁶Pb isotopic systems more appropriate for dating very old samples (like over 1000 Ma). The relatively
- 250 low uncertainties suggest ²⁰⁷Pb/²³⁵U and ²⁰⁷Pb/²⁰⁶Pb isotopic systems are particularly valuable for studying the early history of the earth's crust.

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 escalate with

255 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, for samples older than 3000 Ma, isotopic systems exhibit a significant increase in age uncertainty, often exceeding 500 Ma. While commonly employed strategies such as filtering samples based on acceptable 2σ 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 ²⁰⁷Pb/²⁰⁶Pb isotopic system, where they prove effective in reducing the best age discrepancies (Figure 7i). For instruments, 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 discord ratios over time for different analytical instruments. The SHRIMP method, known for its precision, demonstrates a consistently low discord ratio (Figure 8a). Remarkably, even
- 270 samples with elevated age uncertainties maintain discord ratios below 0.5%, indicating SHRIMP's reliability in dating. LA-ICP-MS suggests an increase in age uncertainty for samples exceeding 1000 Ma but maintains a discord ratio below 0.5% for these samples (Figure 8b). However, a notable disadvantage of LA-ICP-MS is observed for samples within 1000 Ma with low age uncertainties, where a comparatively high discord ratio may exceed 1%, underscoring the need for careful data interpretation in these cases. The ID-TIMS method, while less commonly utilized in sediment dating, exhibits low discord
- 275 ratios (Figure 8c). This suggests that ID-TIMS, despite its limitations, offers robust results for the most precise of dating requirements. Samples analyzed by SIMS appear to exhibit a potential linear relationship between age uncertainty and discord ratio (Figure 8d), providing a straightforward assessment model for data quality in this context. Currently, the LA-ICP-MS method is favored for its lower time consumption in obtaining U-Pb ages. However, the potential for higher discord ratios necessitates additional considerations during data processing to ensure the accuracy and reliability of the results derived from
- this method.
 - The original uncertainty associated with the Lu-Hf dataset predominantly pertains to the analytical outcomes obtained from isotopic measurements. Across all geological periods, the 2σ errors for the isotopic ratios 176 Hf/ 177 Hf, 176 Lu/ 177 Hf, and 176 Yb/ 177 Hf typically fluctuate around 2×10^{-5} , as depicted in Figure S7 in the s Supplement. The measurement ranges for these three isotopes are approximately 2×10^{-2} , indicating a high level of precision in the analytical process (Figure S7 in the
- 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 ϵ Hf(0) and ϵ 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 2σ errors observed in both ϵ Hf(0) and
- 290 EHf(t), as depicted in Figure S8 in the Supplement. These parameters reflect the hafnium isotope composition at the time of



zircon crystallization (ϵ Hf(0)) and at a specific time in the past (ϵ Hf(t)), exhibit 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 ϵ Hf(0) and ϵ Hf(t) are considerably larger than their corresponding 2 σ errors. This discrepancy highlights the precision of the isotopic measurements relative to the calculated uncertainties of the hafnium isotope ratios. The stability of the error over the timescale is particularly noteworthy, suggesting that ϵ Hf(0) and ϵ 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

295

4.3 Spatial and temporal distributions of samples

in geochronological and geochemical analyses.

The spatial distribution biases within the OneDZ database are evident (Fig. 3-4). To delve into the effects of biased distributions, 300 the U-Pb age data was segmented according to geological time sequences and visualized (Fig. 9). Temporal slices reveal that the Qinghai-Tibet Plateau, Alps, Cordillera and Andes mountains are the main sampling areas in the Cenozoic (Fig. 9a-c). In the Mesozoic, mainly sampling regions are similar to places from the Cenozoic (Fig. 9d-f). In Paleozoic, East Asia is obviously over-sampled relative to other regions (Fig. 9g-l). In pre-Cambrian period, East Asia, Europe and Australia contributed major samples (Fig. 9m-n)

- 305 In this study, we also present the first visualization of the temporal distributions of uranium, thorium, and lead concentrations in detrital zircons (Fig. 10). The concentrations of these elements exhibit 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. The data show periodicity, with lead concentration indicating a long-term decline. Notably, there are differences in the estimation of temporal distributions of element concentrations when using Bootstrap and Monte Carlo methods (Fig. 10). Furthermore, beyond elemental
- 310 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 span, the Th/U is larger than 0.4, proving magmatic zircon dominants the detrital zircon (Figure S9 in the Supplement).
- For Lu-Hf isotopes, the ¹⁷⁶Hf/¹⁷⁷Hf isotope keep descending with the ¹⁷⁶Lu/¹⁷⁷Hf and ¹⁷⁶Yb/¹⁷⁷Hf isotopes showing periodic fluctuations (Figure 11). The ϵ Hf(0) displaying continuous fluctuation decline and ϵ Hf(t) periodic fluctuates (Figure 12).

5 Discussion

320

5.1 Evaluate the paleo-globality

Despite the OneDZ database complies 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 lead to overrepresentations of regional samples.





For instance, the spatial analysis of global zircon oxygen isotope has proved that the temporal anomalies in zircon oxygen isotopes were predominantly attributed to the regional samples' imbalance (Sundell et al., 2024). To alleviate the regional imbalance in global zircon data analysis, Puetz et al. (2024) proposed a global representativeness by the ratio of the grid cells activated by gridded data to the total global grid cells. However, this evaluation approach is predicated on the present-day distribution of land and sea. In geological time scale, current geographical pattern does not accurately reflect the samples'

- 325 distribution of land and sea. In geological time scale, 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
- 330 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 time goes on, data tends to concentrate on one or several ancient tectonic plates. Therefore, the evaluation results based on OneDZ, the world's largest detrital zircon database, indicate that the global scope of zircon big data research needs further assessment.
- Following the new paleo-globality evaluating methods, the temporal globality of OneDZ detrital zircon U-Pb data was visualized in Figure 14. The visualization in Figures 14c-e demonstrates that the U-Pb data has achieved spatial coverage across paleo-continents. In calculations, a notable rise (14%) in paleo-globality 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. Estimation at an excessively large scale loses a significant amount of spatial details. Similarly, small-scale grid estimation 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. Although the large volume of the OneDZ,

345 the calculated paleo-globality does not represent the global feature in most geological time (paleo-globality does not equal to 100 percent). For instance, the reconstructed sample distribution in 250 Ma (Figure 13g-h) has a visualized globality. However, the calculated paleo-globality is about 30%~60%. Thus, we suggest that regional data should be treated more carefully when discuss global event. Especially using spatial kernel density evaluation methods is necessary.

5.2 Compare the resampling methods

350 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. 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. The impact of data sparsity is controlled by the 2 σ error. 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 7d-f) and Bootstrap methods (Figure 7g-i), the overall trend of zircon Best age data is consistent. Even on sparse 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
- 365 between the two methods in characterizing the details of time series. In the ²⁰⁶Pb/²³⁸U isotope system, four periodic fluctuations were observed in the Monte Carlo resampling results over the time period of 0-10000 Ma (Figure 7d). The Bootstrap method only shows a slight increase around 500 Ma on the same time scale (Figure 7g). The rest of the time scales show a slow increase. In the ²⁰⁷Pb/²³⁵U 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 7e). The Bootstrap method showed a significant
- 370 decrease around 1500 Ma on the same time scale (Figure 7h). In the ²⁰⁷Pb/²⁰⁶U isotope system, the Monte Carlo resampling results showed a significant decrease around 1500Ma (Figure 7f). On the contrary, the Bootstrap method exhibits a slight periodic decrease (Figure 7i). Although Figure 7 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,
- 375 we compared the results of time resampling methods for zircon U-Pb and Lu Hf system time series data in OneDZ. In zircon U-Pb data, the Bootstrap method shows stronger temporal stability (Figure 10a-c). 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 εHf(0) and εHf(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 Yb/¹⁷⁷Hf 2 σ error time series, due to the significant increase in data volume, the significant difference in the results of resampling methods is relatively small (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 10-12, S7-S7). Considering that Monte Carlo methods typically assume that data follows a normal distribution or is uniform within a time window, data
- 385 typically does not follow a normal distribution within time windows with significant changes in data density. Therefore, the Monte Carlo method is not suitable for time resampling in zircon big data analysis.





Spatial over-sampling introduces another potential bias that has gained attention in the field. 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 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 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. These results suggest that the sampling differences in
- 400 different regions have significantly disappeared (Figure 15c). Direct spatial resampling methods may not fully address spatial imbalances, data enhancements and grid method can substantially reduce spatial biases.

5.3 Implications for database construction and future developments

The construction of the OneDZ database, which employs a crowdfunding mode, has the potential to significantly broaden data coverage. However, crowdfunding mode introduces challenges, such as inconsistencies in data formatting and the risk of

- 405 human errors. To address these issues, a series of Python and MySQL 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 construction process, crowdfunding mode with automatic data cleaning by Python and MySQL code snippets is feasible and greatly improved the efficiency of database construction.
- Moreover, AI tools have played a pivotal role in the data collection and extraction process. Unlike traditional web crawlers, 410 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 DeepShoval 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.
- 415 Computer vision-based AI models show great promise in reducing labor costs and increasing efficiency. Furthermore, the limitations of traditional tools like Excel in storing and managing large datasets, especially given the .xlsx format's known data storage caps and issues with formatting errors, have become increasingly apparent. The need for rapid data retrieval in super large databases also renders Excel inefficient for this purpose. In contrast, MySQL offers both unlimited



420

storage capacity and extremely low retrieval latency, making it a superior tool for constructing super large databases in Earth sciences.

In the OneDZ database, users can quickly retrieve detrital zircon data based on regional search criteria. Regional retrieval is quickly completed through latitude and longitude ranges. This has been achievable in previous datasets. However, OneDZ incorporates information on lithology, stratigraphy and depositing age, which would greatly expand the scope of detrital zircon data. Meanwhile, in conjunction with other Earth science datasets, it is possible to further explore the temporal and spatial avolution patterns of the Earth.

425 evolution patterns of the Earth.

The construction of super large databases in Earth science, utilizing AI and automated scripts, is still in the experimental phase. However, with the rapid advancement of AI, particularly large language models like GPT, we anticipate the development of even more capable models to handle such text-intensive tasks. Additionally, establishing an automated processing platform and designing a user-friendly graphical interface for experts involved in the crowdfunding construction could further reduce

430 labor costs. Progress is being made in developing such a platform and interface, which will be instrumental in advancing the field.

Data availability

The database (Li and Hu, 2025) is freely available via Zenodo at https://doi.org/10.5281/zenodo.14957581. All code snippets in this research are accessible via https://github.com/KeranLi/Global-Detrital-Zircon. The OneDZ web platform is accessible via https://dedc.geoscience.cn/onedz/.

Conclusion

435

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 provides a comprehensive collection of data, encompassing stratigraphy, sedimentary age, isotope geochemical data, and various dating instruments, such as LA-ICP-MS,

440 encompassing stratigraphy, s SHRIMP, SIMS, and TIMS.

Based on this database, we have characterized the errors associated with zircon dating, compared the efficacy of different dating instruments, 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,

445 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 spatiotemporal data resampling methods, Bootstrap methods and SMOTE data augmentation methods may be more suitable.



The construction experience of OneDZ shows that the construction mode of crowdfunding and automated code cleaning is the key to quickly completing database construction. Combined with AI tools and MySQL, the construction and use of databases will be more convenient.

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, 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

455

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

Acknowledgements

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 paleolocations 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.
- 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.



475

485

- 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, 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 2016 international conference on management of data, pp. 2201–2206, https://doi.org/10.1145/2882903.2912574, 2016.
- 480 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, 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.

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.
 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, 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
- 500 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 235U and 238U, Physical Review C, 4, 1889-1906, https://doi.org/10.1103/PhysRevC.4.1889, 1971.
- 505 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.



515

525

540

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

- 510 determinations, Geochimica et cosmochimica acta, 37, 485-494, https://doi.org/10.1016/0016-7037(73)90213-5, 1973.
 - Li, K. and Hu, X.: OneDZ: A Global Detrital Zircon Database and Implications for Constructing Giant Geoscience Database [Data set]. Zenodo. <u>https://doi.org/10.5281/zenodo.15522949</u>, 2025.
 - 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.
 - 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, <u>https://doi.org/10.1038/s41597-022-01730-7</u>, 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.
 - McKenzie, N.R., Horton, B.K., Loomis, S.E., Stocklli, 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.
 - 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, 2243–2261, 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 535 Mineralogy and Petrology, 75: 263–267, https://doi.org/10.1007/BF01166766, 1981.

- 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, 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.



545

550

560

- 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, δ18O zircon, and Sm-Nd whole rock global databases, Scientific Data, 11, 56, https://doi.org/10.1038/s41597-023-02902-9, 2024b.
- Söderlund, U., Patchett, P.J., Vervoort, J.D., Isachsen, C.E.: The 176Lu decay constant determined by Lu–Hf and U–Pb isotope systematics of Precambrian mafic intrusions, Earth Planet. Sci. Lett. 219 (3-4), 311-324, <u>https://doi.org/10.1016/S0012-821X(04)00012-3</u>. 2024.
- 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, https://commons.library.stonybrook.edu/geo-articles/12/, 1985.
- 555 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.
 - 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.
 - 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
- 565 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,
- 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.







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).





580

Figure 2: Temporal distributions of U-Pb and Lu-Hf isotopic records. (a) Kernel density estimate map of U-Pb records (the spatial resolution is 1°×1°); (b) Kernel density estimate map of Lu-Hf records (the spatial resolution is 1°×1°); (c) Era-based distribution of U-Pb samples; (d) Period-based distribution of U-Pb samples; (e) Era-based distribution of Lu-Hf samples; (f) Period-based distribution of Lu-Hf samples.









Figure 3: Visualizations of the spatial, temporal and strata information in U-Pb dataset. (a)-(b) Major geographic/geological 585 description; (c)-(d) Minor geographic/geological description; (e)-(f) Group-Formation-Member records; (g)-(h) Locality/Sedimentary profile.







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.







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.



595 Figure 6: Ages errors of different isotopic systems. (a) ²⁰⁶Pb/²³⁸U; (b) ²⁰⁷Pb/²³⁵U; (c)²⁰⁷Pb/²⁰⁶Pb.







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.







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







Figure 9: Spatial-temporal distribution of U-Pb age data.







605 Figure 10: Discord ratio varying with time by different instruments. (a) SHRIMP; (b) LA-ICP-MS; (c) ID-TIMS; (d) SIMS.



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

610







Figure 12: Temporal variations of ε Hf uncertainties in Lu-Hf dataset. (a) ε Hf(0) with bootstrap resampling; (b) ε Hf(t) with bootstrap resampling; (c) ε Hf(0) with Monte-Carlo resampling; (d) ε Hf(t) with Monte-Carlo resampling.







615 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°×1°). (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.







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







 $\label{eq:starsest} Figure 15: The resampling frequency of different methods. (a) Monte-Carlo method; (b) SMOTE-Monte-Carlo method; (c) 12°\times12° grid-SMOTE-Monte-Carlo method.$



625

Table 1: Construction methods comparison of three typical zircon databases

Dataset	Methods	Inform	nation types	Dat	a cleaning	Adapting A	I tool M	anage data
OneDZ	Crowdfunding	Data o	rigin, Spatial	Ру	thon and	Yes		MySQL
		informatio	on, Stratigraph	ic l	MySQL			
		informa	ation, Isotopic	sc	ripts and			
		inf	formation	8	rtificial			
				с	hecking			
Wu et al., 2024	Directly	Data o	rigin, Spatial	A	Artificial	Not mentio	oned	Excel
	collecting	informa	ation, Isotopic	с	hecking			
		inf	formation					
Puetz et al., 2024	Directly	Data o	rigin, Spatial	A	Artificial	Not mentio	oned	Excel
	collecting	informatio	on, Stratigraph	ic c	hecking			
		informa	ation, Isotopic					
		inf	ormation					
Table 2: Data compar	ison of three typi	cal zircon data	abases					
Dataset	Field	Field	Reference	Sample	Region	Geological	Valid U-	Valid
	number of	number of	(10 ⁴)	(104)		unit	Pb age	Lu-Hf
	U-Pb	Lu-Hf					with GPS	data with
							(10^{6})	GPS
								(10 ⁵)
OneDZ	71	86	5.4	31	215	1348	1.8	2.7
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	2.1

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.

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

Dataset	Field Name	Corresponding field in	Corresponding field in Wu et	Proportion
		Puetz et al., 2024	al., 2024	
U-Pb	Lead_Author	Lead_Author	Author_surname and	100.00%
			Author_given_name	
	Year	Year	Year_publication	100.00%
	Journal	Journal	Journal	100.00%





	Vol.	Vol.	Volume	96.99%
	Pages	Pages	First_page and Last_page	95.93%
	Title	Title	Title	100.00%
	Web_Link	Web_link	DOI	54.50%
Lu-Hf	Lead_Author	Lead_Author	Lead_Author	100.00%
	Year	Year	Year	100.00%
	Journal Journal Vol. Vol.	Journal	100.00%	
		Vol.	Volume	96.48%
	Pages	Pages	Pages / Article No.	96.48%
	Title	Title	Title	100.00%
	Web_Link	Web_Link	Web Link	100.00%

630 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)

		Corresponding field in Puetz et	Corresponding field in Wu et		
Dataset	Parameter	corresponding nero in ruetz et	Conceptioning netu in will et	Proportion	
		al., 2024 al., 2024			
U-Pb	Published_Sample_ID	Smaple_ID	Sample_number	26.88%	
	Country_State	Country/Small Region		26.91%	
	Region	Large Region		79.75%	
	Continent	Continent		75.96%	
	Maine Carlada Davidia	Major Geographic-Geologic		26.000/	
	Major_Geologic_Description	Description		26.90%	
	Minor_Geologic_Description	Minor Geologic-Geographic Unit		21.45%	
	Group			9.97%	
	Formation			20.94%	
	Member			11.30%	
	Locality	Locality		66.10%	
	Sedimentary profile			4.61%	
	Latitude	Latitude	Latitude	96.23%	
	Longitude	Longitude	Longitude	96.23%	
	Depos. Age (Period)			11.31%	
	Depos. Age (Epoch)			10.95%	
	Depos. Age (Age)			10.14%	



	Max. Depos. Age (Ma)	Max. Stratigr. Age (Ma) (detrital only)		10.15%
	Est. Depos. Age (Ma)	Est. Stratigr. Age (Ma) (detrital only)		11.51%
	Min. Depos. Age (Ma)	Min. Stratigr. Age (Ma) (detrital only)		9.14%
Lu-Hf	Published_Sample_ID	Ref. No.	Published Sample_ID	21.47%
	Country_State		Country/Small Region	75.85%
	Region			78.76%
	Continent		Continent	73.99%
	Major_Geologic_Description		Major_Geologic_Description	68.78%
	Minor_Geologic_Description		Minor_Geologic_Description	67.01%
	Group			5.05%
	Formation			10.05%
	Member			0.80%
	Locality		Locality	80.86%
	Sedimentary profile			1.82%
	Latitude		Latitude	99.03%
	Longitude		Longitude	99.03%
	Depos. Age (Period)			12.28%
	Depos. Age (Epoch)			8.09%
	Depos. Age (Age)			3.92%
	Max. Depos. Age (Ma)			8.85%
	Est. Depos. Age (Ma)		Est. Strat. Age (Ma)	46.17%
	Min. Depos. Age (Ma)			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	Corresponding field in	Corresponding field in Wu	Proportion
	Puetz et al., 2024	et al., 2024	
Mass_Spectrometer	Mass Spectrometer	Instrument	57.02%
Spectrometer_Location	Spectrometer Location		58.42%
Institution	Institution		61.09%
U_Pb_Record_Count	Accepted records		60.97%



Grain_ID	Sample&Grain			98.16%
Spot	Spot			75.95%
Spot_diam	Spot diam. (μm)		12.14%
206Pb_238U_isotope_ratio	206Pb/238U	ratio	isotope206Pb/238U	69.16%
206Pb_238U_isotope_uncertainty_1sigma	206Pb/238U	1σ	isotope206Pb/238U_σ	63.26%
	uncert			
206Pb_238U_isotope_uncertainty_2sigma				63.26%
207Pb_235U_isotope_ratio	calculate	d	isotope207Pb/235U	21.31%
	207Pb/235U	ratio		
207Pb_235U_isotope_uncertainty_1sigma	207Pb/235U	1σ	isotope207Pb/235U_σ	17.06%
	uncert			
207Pb_235U_isotope_uncertainty_2sigma				17.06%
207Pb_206Pb_isotope_ratio	207Pb/206Pb	ratio	isotope207Pb/206Pb	6.77%
207Pb_206Pb_isotope_uncertainty_1sigma	207Pb/206	Pb	isotope207Pb/206Pb_o	5.00%
	1σ uncer	t		
207Pb_206Pb_isotope_uncertainty_2sigma			Optional	5.00%
208Pb_232Th_isotope_ratio	isotope208Pb/	232Th	Optional	3.09%
208Pb_232Th_isotope_uncertainty_1sigma	isotope208Pb/2	232Th_	Optional	3.09%
	σ			
208Pb_232Th_isotope_uncertainty_2sigma			Optional	3.09%
Published_206Pb_238U_age	206Pb/238U	age	age206Pb/238U	97.94%
	(Ma)			
Published_206Pb_238U_age_uncertainty_1sigma			age206Pb/238U_σ	90.93%
Published_206Pb_238U_age_uncertainty_2sigma	206Pb/238U	2σ		90.93%
	uncert			
Published_207Pb_235U_age	207Pb/235U	age	age207Pb/235U	90.94%
	(Ma)			
Published_207Pb_235U_age_uncertainty_1sigma			age207Pb/235U_σ	90.94%
Published_207Pb_235U_age_uncertainty_2sigma	207Pb/235U	2σ	Optional	90.94%
	uncert			
Published_207Pb_206Pb_age	207Pb/206Pb	age	age207Pb/206Pb	93.38%
	(Ma)			

35



Published_207Pb_206Pb_age_uncertainty_1sigm		age207Pb/206Pb_σ	93.38%
a			
Published_207Pb_206Pb_age_uncertainty_2sigm	207Pb/206Pb		93.38%
a	2σ uncert		
Best_Age			99.99%
Best_Age_uncertainty_1sigma			95.56%
Best_Age_uncertainty_2sigma			95.56%
Discord_ratio			25.12%
U_ppm			21.56%
Th_ppm			21.56%
Pb_ppm			21.56%

Table 6: Data specifications of the Lu-Hf isotopic system (the proportion was calculated by number of valid items divided the635number of total items)

Field Name	Corresponding field in Puetz et al., 2024	Proportion
Mass_Spectrometer		10.26%
Spectrometer_Location		14.28%
Institution		3.83%
Lu_Hf_Record_Count		9.57%
Grain_ID		20.66%
Spot		5.77%
Spot_diam		3.87%
U_Pb_Age	U-Pb Age (Ma)	99.86%
U_Pb_Age_uncertainty_1sigma		13.02%
U_Pb_Age_uncertainty_2sigma		13.02%
176Hf_177Hf	176Yb/177Hf sample ratio	99.67%
176Hf_177Hf_uncertainty_1sigma		13.82%
176Hf_177Hf_uncertainty_2sigma	176Yb/177Hf 2σ	13.82%
176Lu_177Hf	176Lu/177Hf sample ratio	99.86%
176Lu_177Hf_uncertainty_1sigma		55.40%
176Lu_177Hf_uncertainty_2sigma	176Lu/177Hf 2σ	55.40%
176Yb_177Hf	176Hf/177Hf sample ratio	76.34%
176Yb_177Hf_uncertainty_1sigma		41.67%
176Yb_177Hf_uncertainty_2sigma	176Hf/177Hf 2σ	41.67%





178Hf_177Hf		4.72%
178Hf_177Hf_uncertainty_1sigma		0.23%
178Hf_177Hf_uncertainty_2sigma		0.23%
180Hf_177Hf		75.24%
180Hf_177Hf_uncertainty_1sigma		75.24%
180Hf_177Hf_uncertainty_2sigma		75.24%
176Hf_177Hf_initial		9.75%
176Hf_177Hf_initial_uncertainty_1sigma		0.24%
176Hf_177Hf_initial_uncertainty_2sigma		0.24%
εHf(0)		13.77%
εHf(0)_uncertainty_1sigma		4.26%
εHf(0)_uncertainty_2sigma		4.26%
ɛHf(t)	$\epsilon H f(t)$ calc	99.90%
ϵ Hf(t)_uncertainty_1sigma		86.54%
eHf(t)_uncertainty_2sigma	$\epsilon H f(t) 2\sigma$ calc	86.54%
TDM1	TDM1 (Ma) calc	95.49%
TDM1_uncertainty_1sigma		8.42%
TDM1_uncertainty_2sigma		8.42%
TDM2	TDM2 (Ma) calc	95.31%
TDM2_uncertainty_1sigma		4.54%
TDM2_uncertainty_2sigma		4.54%
176Hf_177Hf_Chur		6.21%
176Hf_177Hf_DM		4.10%
176Lu_177Hf_Chur		0.28%
176Lu_177Hf_DM		2.15%