



26 **1 Introduction**

27 Geoanalytical data includes measurements of major and trace elements, rare earth elements (REEs), isotopes, and
28 structures and morphology, of geological samples analyzed by various analytical instruments such as ICP (LA-ICP-MS,
29 ICP-MS), EPMA, SIMS, SEM, TEM, XRF, and XRD. Geoanalytical data effectively reflect the material composition,
30 internal structure, external characteristics, interaction and evolution history of the Earth and represent the most important
31 support for geological researchers in their aim to understand the Earth and exploit its resources for the survival and
32 development of human society. Enormous financial, material, and human resources have been invested into the
33 geological surveys and geoanalysis required to acquire more comprehensive and abundant geoanalytical data. Over time,
34 tremendous volumes of geoanalytical data have been created, and these volumes continue to increase at a high rate. It is
35 of paramount importance that this data is curated effectively and that adequate background information, such as sample
36 description, sampling information, and analysis information, is included, so that geological researchers can utilize the
37 data according to their requirements. This will also facilitate the reutilization of the precious geoanalytical data. In
38 addition, with the accumulation of large volumes of data, statistical analysis and data mining can be conducted on these
39 data to provide a more comprehensive and advanced scientific understanding of the Earth. Hence, a variety of
40 geoanalytical databases aimed at managing, sharing and reutilizing geoanalytical information have been constructed and
41 are used as advanced tools in geological studies. The analysis and comparison of existing geoanalytical data models, as
42 well as the development of improved models, is therefore a worthy and significant study to be conducted.

43 Over the last decades, several studies of geoanalytical data models have been conducted. As early as 1977, George
44 Van Trump and colleagues described a data model for environmental geochemical surveying and mineral resource
45 exploration in the United States of America (Jr and Miesch, 1977). Lehnert et al. (2000) suggested a data model for the
46 storage of global geochemical data of rocks. Their data model provides a complete summary of essential geochemical
47 data contents and a robust structure with relational database management system (RDBMS). Numerous databases such as
48 GEOROC, NAVDAT, and PetDB have since been constructed based on this model, and it is used by geological
49 researches worldwide. In particular PetDB has been used for a considerable amount of high-impact research such as
50 *Nature* (Brandl et al., 2013; Carbotte et al., 2013; Cheng et al., 2016; Dick and Zhou, 2014; Helo et al., 2011; Hoernle et
51 al., 2011; Kamenov et al., 2011; Kelley, 2014; Samuel and King, 2014; Schlindwein and Schmid, 2016; Straub et al.,
52 2009) and *Science* (Cottrell and Kelley, 2013; Greber et al., 2017; Joy et al., 2012; Kelley and Cottrell, 2009; McNutt et
53 al., 2016). A limitation of existing geoanalytical data models is their specificity to particular applications or geological



54 domains and their focus on the description and curation of only a certain portion of geoanalytical data. For example,
55 RU_CAGeochem is specifically focused on major and trace element concentrations and Sr, Nd, and Pb isotopic ratios of
56 American volcanic rocks (Carr et al., 2014). Another database is focused on lead-isotopes of copper ores from the
57 south-eastern Alps (Artioli et al., 2016). Many other examples of similarly specific geoanalytical databases and
58 associated models exist (e.g. Artioli et al., 2016; Hellström, 2016; Lopes et al., 2014; Siegel et al., 2012; Strong et al.,
59 2016). The consequence of this development is that each database exists as a separate island and it is difficult for
60 researchers to communicate and integrate geoanalytical data between databases. In particular, every time a database is
61 constructed, a data model has to be redesigned. This consumes considerable amounts of time and prolongs the
62 development cycle. In addition, the vast majority of models are designed based on relational models, which focus on the
63 construction of relations between different data categories. When users query and utilize the geoanalytical data from
64 different dimensions, these types of models utilize complicated joints between different tables to query the target data,
65 which decreases efficiency as the amount of data increases. However, the exploration of such data models including the
66 background items, have laid a solid foundation for later study of advanced geoanalytical data models.

67 At present, the development of various new techniques provides us with the opportunities to design more
68 comprehensive and advanced geoanalytical data models. In this study, we introduce a novel, universal and efficient
69 geoanalytical data model. First, we provide an overview of geoanalytical methods and applications to summarize the
70 geoanalytical data available. Then, we design universal data attributes based on this data and develop a
71 multi-dimensional data model. Finally, we evaluate the model to validate its efficiency.

72 **2 Overview of geoanalytical data contents**

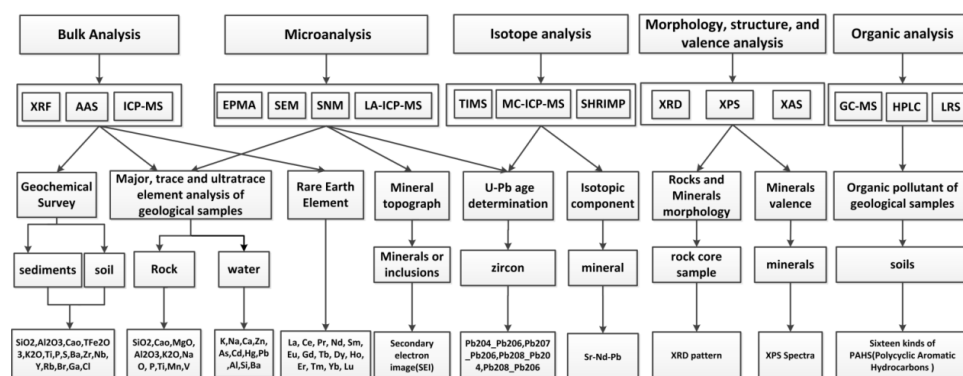
73 In recent years, many new geoanalytical methods and instruments have been developed, creating novel kinds of
74 data (Linge et al., 2017). A truly universal data model should have the ability to accommodate all kinds of
75 geoanalytical data. In addition, the data model should be capable of making all stored data readily available for
76 reutilization by geological researchers. In order to develop a model with such capabilities, a comprehensive set of
77 geoanalytical data, together with related background information required for reutilization of the data, was
78 summarized and categorized, as outlined below.

79 First, analytical techniques and their applications were studied to comprehensively summarize geoanalytical
80 measurement data. This process is outlined in Figure 1. Because of the great diversity of analytical methods and
81 geological applications, Figure 1 only shows a few examples to indicate the method adopted in this paper. The five



82 categories, namely, bulk analysis; microanalysis; isotope analysis; morphology, structure, and valence analysis; and
 83 organic analysis, were divided according to the analytical technique used. In this way, data from each category was
 84 categorized according to analytical instruments (e.g., SEM, SNM, and EPMA for microanalysis). In the next step,
 85 the data were grouped according to geological applications. The comprehensive list of geoanalytical measurement
 86 data items used in the present study, compiled from a thorough literature review, is presented in Figure 2. In the
 87 case of bulk analysis, most measurements ultimately provided major, trace and ultra-trace element concentration
 88 data. Microanalysis can yield data of elemental concentrations in a micro-region, as well as structural information
 89 of geological samples acquired by secondary electron and backscattered electron techniques, commonly stored as
 90 image files. For geochronology and stable isotopic analysis (GSI analysis), most measurement data are isotopic
 91 ratios. For morphology, structure and valence analysis (MSV analysis), the most common measurement data are
 92 image files such as XPS spectra or XRD patterns. Organic analysis is a new analytical method, which is used for
 93 the analysis of environmental geological samples. The most common application of this method in the geological
 94 literature is the analysis of the sixteen kinds of Polycyclic Aromatic Hydrocarbons PAHs in soils.

95



96

97

Figure 1: Process of summarizing the geoanalytical data contents.

98

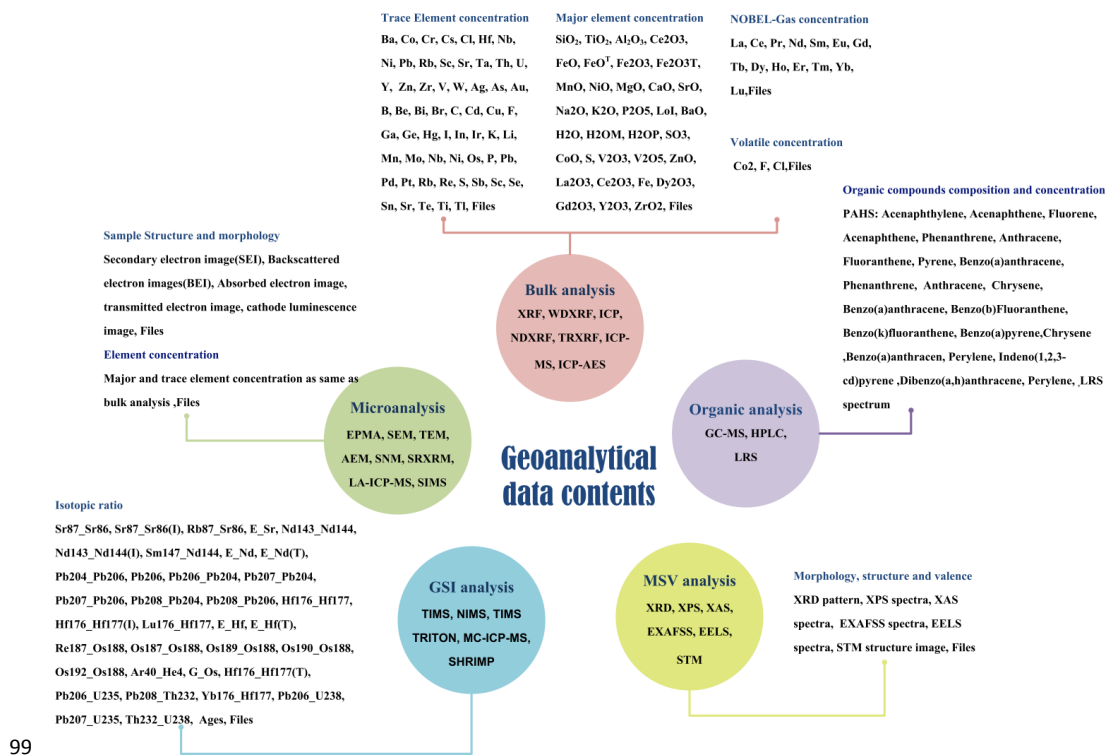


Figure 2 lists of geoanalytical methods and measurement data items

Background information describing the analyzed samples and data quality has to be incorporated, because it is indispensable for proper evaluation, efficient recovery, and sorting of the compiled data. Hence, background metadata are summarized based on the investigations of the geological researchers and the contents of existing databases (Adcock et al., 2003; Lehnert et al., 2000). Table 1 lists details of the background metadata used during the present study. In this study, the Background metadata are divided into three parts: sample metadata provide geological researchers with information about geological materials, sampling metadata provides information about environmental conditions in the field, and quality metadata allows geological researchers to make an assessment of data quality and usability (Table 1). The background metadata items listed in Table 1 are the most essential information required for every kind of geoanalytical measurement data. More specific attributes are not included in our model.

Table 1 Background metadata of geoanalytical data

Background Metadata



Sample metadata	sample type, sample name, sample description
Sampling metadata	sampling site description, longitude, latitude, sampling methods, sampling depth, sampling institution
Quality metadata	laboratory name, project name, publish link, analytical method, analysts

112 **Geoanalytical data modelling**

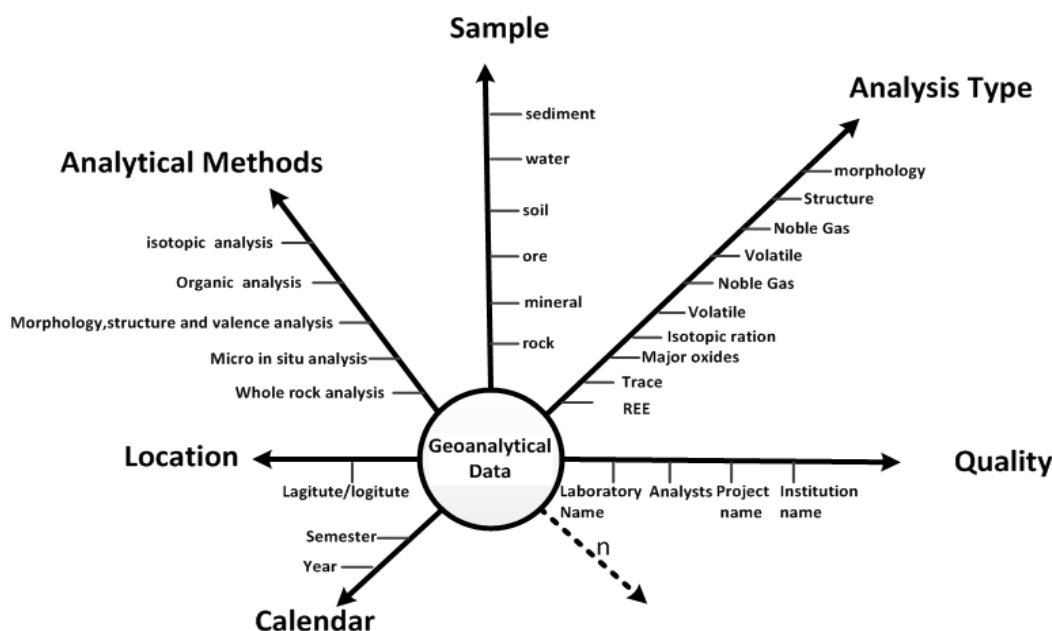
113 This section outlines how the novel geoanalytical data model was designed, utilizing the data summarized above.
114 Despite their limitations, currently relational data mode is the most commonly used pattern for geoanalytical data
115 models. Relational data mode construct relations between each group of data within the database. This means that
116 more data categories inevitably lead to much more data relations, increasing storage demands and the time required
117 to query the database. Compared to such conventional relational data models (Beynon-Davies, 2004),
118 multi-dimensional models (MDM), which are widely utilized during the development of big data science and data
119 mining, are single subject-oriented sources for analyzing data based on various dimensions (Niemi and Hirvonen,
120 2003). Multi-dimensional modelling approaches share characteristics with Fast Analysis of Shared
121 Multidimensional Information (FASMI). In particular, MDM offers the advantage of a relatively simple and
122 straight-forward database design, which nevertheless supports powerful analyses, and is relatively well understood
123 by the end users (Hoberman, 2005). As a modelling framework, MDM has a conceptual and a logical phase of
124 design, composed of a fact table and several dimension tables (Höpken et al., 2013). Facts comprise numeric and
125 additive characteristics of the data, which can be accumulated along multiple dimensions. Frequently, researchers
126 are interested in analyzing geoanalytical measurement data from different metadata perspectives. Hence, the MDM
127 approach is ideally suited for the design of geoanalytical data models. Here, the geoanalytical data are the fact data,
128 and other background information are dimensions data. The use of the MDM modelling framework applied in the
129 present study will allow geological researchers to rapidly analyze geoanalytical data based on numerous meta-data
130 criteria.

131 **2.1 Conceptual data model (CDM)**

132 A conceptual data model (CDM) includes the definition of its universal attributes and a rough design of its structure.
133 It represents the primary phase of data model design, independent from the detailed techniques of computer
134 systems. Figure 3 presents the multi-dimensional CDM we developed for geoanalytical data. Here, the abstraction



135 of universal concepts present in geoanalytical data, the model becomes more flexible and universally applicable.
 136 The geoanalytical data are placed in the center of the model, in the form of a facts table. The associated background
 137 information is categorized and abstracted as various dimensions which are represented by different axes in Figure 3.
 138 The six dimensions of our CDM are sample, analysis type, analytical methods, location, time and quality. This
 139 arrangement allows geological researchers to analyze geoanalytical data from six different dimensions or any
 140 combination thereof. The marks in each dimension represent the detailed measurement conditions. The “n”
 141 dimension is an expansible dimension, which can be added according to the specific model application.



142

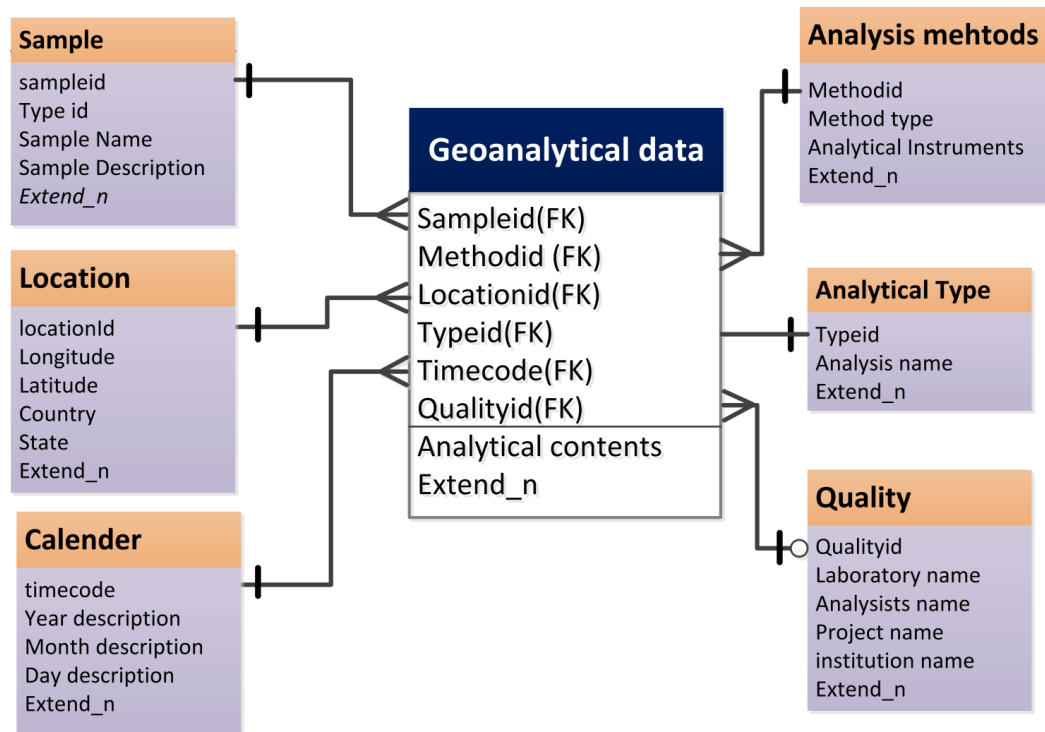
143 **Figure 3: Multi-dimensional CDM for geoanalytical data.**

144 **2.2 Logical data model (LDM)**

145 A logical data model (LDM) is a CDM written in unified modelling language (UML) (Evans et al., 2014). Logical
 146 model design leads to a logical scheme, defining objects, attributes, and relationships (Chmura and Heumann,
 147 2005). The LDM scheme can be easily implemented by any DBMS. Figure 4 shows the LDM scheme designed for
 148 geoanalytical data. Each box in the LDM represents an object, and items in the box are its attributes. The relations
 149 between object are represented with lines. There are three kinds of symbols associated with the lines. The short line



150 denotes “1”, the circle denotes “0” which means “maybe”, and the triangle denotes “many”. Lines and symbols
 151 define the relations between objects. The additional notation foreign key (FK) is added if the attribute in one object
 152 uniquely identifies an attribute in another object. For example, the sample ID in the geoanalytical data object is a
 153 foreign key of Sample_id in the sample object, because they have the same value. By means of this foreign key, the
 154 data contents of the two objects are connected. For each object, a few extended attributes are added (*Extend_n* in
 155 Figure 4). This feature allows developers to add database specific attributes to this model, increasing its flexibility
 156 and universal applicability.



157
 158

Figure 4: LDM of the geoanalytical data model.

159 3 Implementation and Evaluation

160 In order to evaluate the performance of our model, we carried out a comparison experiment with the widely used
 161 Lehnert rock geochemical data model (Lehnert et al., 2000). In order to conduct the experiment, a physical data
 162 model (PDM) needed to be created with a database management system. As RDBMS is the most common
 163 techniques used in geoanalytical databases, MySQL which is a widely used RDBMS was adopted to implement the



164 two models. A specific data item (rocktype: andesite; location: sycamore Hall; latitude: 36.27.12N, longitude:
 165 83.34.12W; Institution: Jilin University; Method: ICP-MS; SiO2:58.9; FiO2:1.13) was used as test data and tables
 166 related to the data contents were implemented. We analyzed the two models from two perspectives: developers and
 167 users. For developer, the comparison of the PDM structure is shown in Figure 5, and query operation descriptions
 168 are presented Figure 6. The comparison clearly indicates that the geoanalytical data model is more succinct than
 169 rock data model, and saving time and computer resources. Three model performance indicators (insert time, storage
 170 space usage, and retrieval time) were evaluated with the increasing of amounts of data. The results are shown in
 171 Figures 7, 8, and 9, respectively. Figure 7 shows clearly that the process of data insertion is considerably faster for
 172 the geoanalytical data model, when compared to the rock data model. Figure 9 shows clearly that the storage space
 173 usage is relatively less than rock data model. In the case of data query (Figure 8), the difference in time
 174 consumption is even more striking. With increasing amount of data items, the query time of the geoanalytical data
 175 model remains very fast and efficient. In contrast, for the rock data model query time costs increased exponentially
 176 with the increasing amount of data items.

177
 178
 179
 180
 181

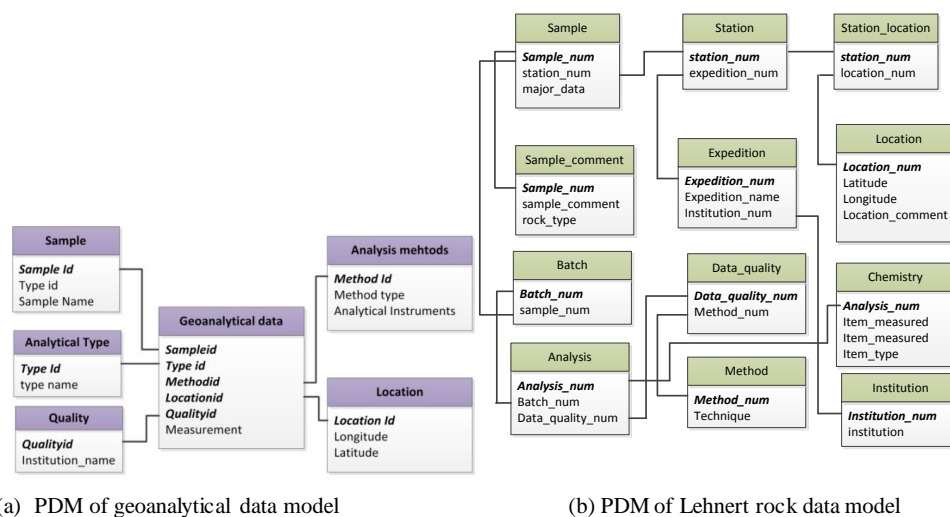


Figure 5: PDM structure comparison with the Lehnert rock data model.



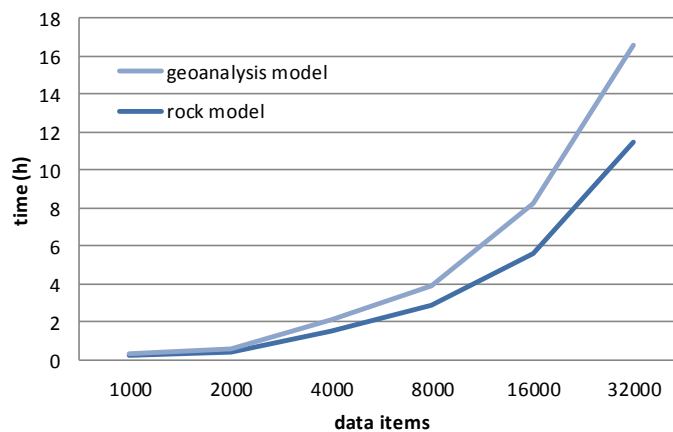
SQL comparison	
Insert operation	
Geoanalysis model	Rock model
<ul style="list-style-type: none"> • insert into sample(sample_num,station_num) values(10, 50) • insert into analysis(analysis_num,batch_num,data_quality_num) values(30,10,90) • insert into batch(batch_num,sample_num) values(20, 10) • insert into station(station_num) values(50) • insert into institution(institution_num,institution_name) values(70,jilin university) • insert into station_location(station_num,location_num) values(50, 80) • insert into location_num(station_num, longitude, latitude) values(30,83.34.10,36.27.12) • insert into data_quality(data_quality_num,method_num) values(90, 40) • insert into method(method_num,technique) values(40, ICPMS) • insert into chemistry(analysis_num,,item_measured,item_type) values(30,,Sio2,59) 	<ul style="list-style-type: none"> • insert into sample (sampleid, sample_name) values(10,sulphate) • insert into quality(qualityid,institution_name) values(20, Jilin university) • insert into location(locationid, longitude, latitude) values(30,83.34.10,36.27.12) • insert into analysis_methods(methodid,analytical_instruments) values(40, ICPMS) • insert into measurement(sample_id,methodid,Locationid,qualityid,item,value) values(10,40,30,20,Sio2,59)
Query operation	
Geoanalysis model	Rock model
<pre>select l.location_num, l.latitude, l.longitude, s.station_num, s.location_num, t.station_num, a.sample_num, a.station_num, b.batch_num, b.sample_num, y.batch_num, y.analysis_num, c.analysis_num, c.item_measured, c.item_type from location l, station_location s, station t, sample a, batch b, analysis y, chemistry c where l.latitude= '36.27.12' and l.longitude='83.34.10' and l.location_num=s.location_num and s.station_num=t.station_num and t.station_num=a.station_num and a.sample_num=b.sample_num and b.batch_num=y.batch_num and y.analysis_num=c.analysis_num;</pre>	<pre>select l.locationid, l.longitude, l.latitude, m.locationid, m.item, m.value from location l, measurement m where l.latitude= '83.34.10' and l.longitude='36.27.12' and l.locationid=m.locationid;</pre>

182

183

184

Figure 6: Comparison of insert and query operations in SQLs.



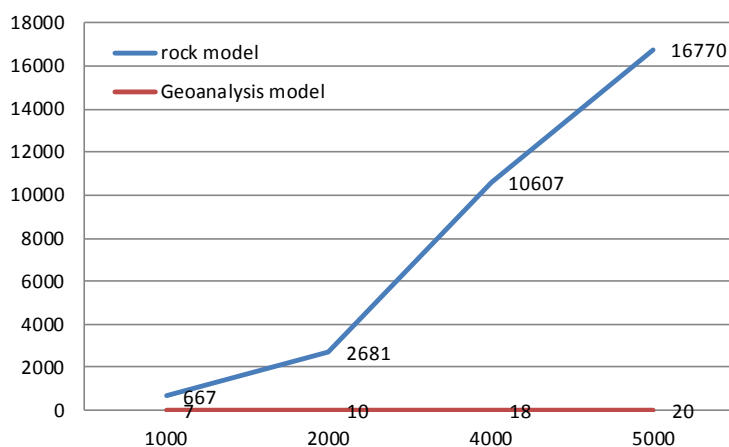
185

186

187

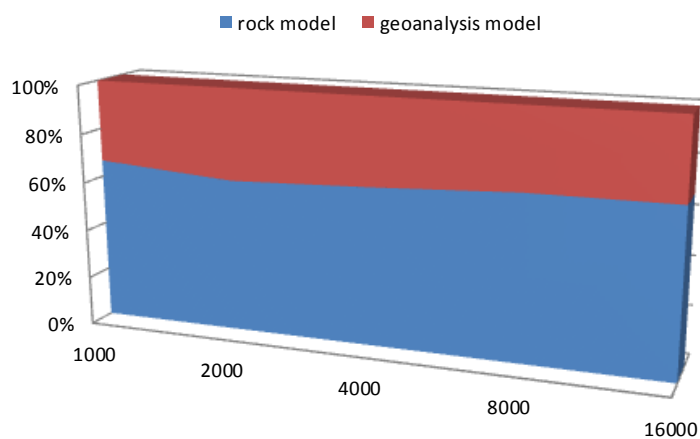
188

Figure 7: Time spent on data insert operations with increasing amounts of data.



189
 190
 191
 192

Figure 8: Time requirement for data queries (latitude and longitude).



193
 194
 195

Figure 9: Space usages of the two data models with increasing amount of data items.

196 4 Conclusions

197 The geoanalytical data model presented herein is flexible and appropriate for a broad range of applications to
 198 geoanalytical data. The model has the following general characteristics:

199 [1] Its universality allows the model to accommodate any type of geoanalytical data for various geological



200 materials, as well as all significant metadata.

201 [2] The adoption of a multi-dimensional data model framework provides geological researchers with the ability to
202 analyze geoanalytical data from different dimensions. In addition to the sample description and location criteria
203 commonly used in existed databases, this model provides four additional query criteria (method, quality, time, and
204 analysis).

205 [3] There are minimum data relations between different objects. Relations between different background metadata
206 objects have been avoided in order to construct robust relations between background metadata and measurement
207 data. This increases the model's efficiency when geoanalytical data are inserted or queried while simultaneously
208 decreasing its space usage.

209

210 It is hoped that the design of this model will allow for the unified construction of geoanalytical databases. The
211 model enables the accumulation and integration of significant amounts of diverse geoanalytical data. By utilization
212 of the big data analysis techniques described in our study, geological researchers could analyze geoanalytical data
213 with high efficiency and develop novel methods to conduct Earth science studies.

214 **Acknowledgments:**

215 This work was financially supported by National Major Scientific Instruments and Equipment Development
216 Special Funds [2016YFF0103303, 2011YQ050069] and CGS research fund [JYYWF20181702].

217 **Reference**

218 Adcock, S.W., Grunsky, E. and Laframboise, R., 2003. A Universal Geochemical Survey Data Model.

219 Beynon-Davies, P., 2004. Relational Data Model.

220 Brandl, P.A., Regelous, M., Beier, C. and Haase, K.M., 2013. High mantle temperatures following rifting caused by
221 continental insulation. *Nature Geoscience*, 6(5): 391-394.

222 Carbotte, S.M., Marjanović, M., Carton, H., Mutter, J.C., Canales, J.P., Nedimović, M.R., Han, S. and Perfit, M.R., 2013.
223 Fine-scale segmentation of the crustal magma reservoir beneath the East Pacific Rise. *Nature Geoscience*, 6(10):
224 866-870.

225 Cheng, H., Zhou, H., Yang, Q., Zhang, L., Ji, F. and Henry, D., 2016. Jurassic zircons from the Southwest Indian Ridge.



- 226 Sci Rep, 6: 26260.
- 227 Chmura, A. and Heumann, J.M., 2005. Logical Data Modeling. *Integrated*, 5(2): 179-203.
- 228 Cottrell, E. and Kelley, K.A., 2013. Redox heterogeneity in mid-ocean ridge basalts as a function of mantle source.
- 229 *Science*, 340(6138): 1314.
- 230 Dick, H.J.B. and Zhou, H., 2014. Ocean rises are products of variable mantle composition, temperature and focused
- 231 melting. *Nature Geoscience*, 8(1): 68-74.
- 232 Evans, A., France, R., Lano, K. and Rumpe, B., 2014. The UML as a Formal Modeling Notation. *Computer Standards &*
- 233 *Interfaces*, 19(7): 325-334.
- 234 Greber, N.D., Dauphas, N., Bekker, A., Ptáček, M.P., Bindeman, I.N. and Hofmann, A., 2017. Titanium isotopic evidence
- 235 for felsic crust and plate tectonics 3.5 billion years ago. *Science*, 357(6357): 1271-1274.
- 236 Höpken, W., Fuchs, M., Höll, G., Keil, D. and Lexhagen, M., 2013. Multi-Dimensional Data Modelling for a Tourism
- 237 Destination Data Warehouse.
- 238 Helo, C., Longpré, M.A., Shimizu, N., Clague, D.A. and Stix, J., 2011. Explosive eruptions at mid-ocean ridges driven
- 239 by CO₂-rich magmas. *Nature Geoscience*, 4(4): 260-263.
- 240 Hoberman, S., 2005. Data Modeling Essentials. *Dm Review*, 131(8): 654 - 660.
- 241 Hoernle, K., Hauff, F., Werner, R., Bogaard, P.V.D., Gibbons, A.D., Conrad, S. and Müller, R.D., 2011. Origin of Indian
- 242 Ocean Seamount Province by shallow recycling of continental lithosphere. *Nature Geoscience*, 4(12): 883–887.
- 243 Joy, K.H., Zolensky, M.E., Nagashima, K., Huss, G.R., Ross, D.K., McKay, D.S. and Kring, D.A., 2012. Direct detection
- 244 of projectile relics from the end of the lunar basin-forming epoch. *Science*, 336(6087): 1426.
- 245 Jr, V.T. and Miesch, A.T., 1977. The U.S. geological survey mass-statpac system for management and statistical reduction
- 246 of geochemical data. *Computers & Geosciences*, 3(3): 475-488.
- 247 Kamenov, G.D., Perfit, M.R., Lewis, J.F., Goss, A.R., Jr, R.A. and Shuster, R.D., 2011. Ancient lithospheric source for
- 248 Quaternary lavas in Hispaniola. *Nature Geoscience*, 4(8): 554-557.
- 249 Kelley, K.A., 2014. Inside Earth Runs Hot and Cold. *Science*, 344(6179): 51-52.
- 250 Kelley, K.A. and Cottrell, E., 2009. Water and the oxidation state of subduction zone magmas. *Science*, 325(5940):
- 251 605-7.
- 252 Lehnert, K., Su, Y., Langmuir, C.H., Sarbas, B. and Nohl, U., 2000. A global geochemical database structure for rocks.
- 253 *Geochemistry Geophysics Geosystems*, 1(5): 179-188.
- 254 Linge, K.L., Bédard, L.P., Bugoi, R., Enzweiler, J., Jochum, K.P., Kilian, R., Liu, J., Marin-Carbonne, J., Merchel, S. and



- 255 Munnik, F., 2017. GGR Biennial Critical Review: Analytical Developments Since 2014. *Geostandards &*
256 *Geoanalytical Research*, 36(4): 337-398.
- 257 Mcnutt, M.K., Lehnert, K., Hanson, B. and Nosek, B.A., 2016. Liberating field science samples and data. *Science*,
258 351(6277): 1024.
- 259 Niemi, T. and Hirvonen, L., 2003. *Multidimensional data model and query language for informetrics*. John Wiley & Sons,
260 Inc., 939–951 pp.
- 261 Samuel, H. and King, S.D., 2014. Mixing at mid-ocean ridges controlled by small-scale convection and plate motion.
262 *Nature Geoscience*, 7(8): 602-605.
- 263 Schlindwein, V. and Schmid, F., 2016. Mid-ocean-ridge seismicity reveals extreme types of ocean lithosphere. *Nature*,
264 535(7611).
- 265 Straub, S.M., Goldstein, S.L., Class, C. and Schmidt, A., 2009. Mid-ocean-ridge basalt of Indian type in the northwest
266 Pacific Ocean basin. *Nature Geoscience*, 2(4): 286-289.
- 267