4    # Data-Driven Abductive Discovery in Mineralogy

5    **Robert M. Hazen[1]\***

6    [1]*Geophysical Laboratory, Carnegie Institution of Washington,*

7    *5251 Broad Branch Road NW, Washington, D. C. 20015, USA.*

8

9    **ABSTRACT**

10    Traditional pathways to discovery in Earth sciences rely in large measure on deductive

11    and inductive approaches, by which measurements and observations are made in the

12    context of established principles or testable, predictive hypotheses about the natural

13    world. Vast, but largely untapped, Earth and life science data resources offer a potentially

14    revolutionary alternative "abductive" approach to investigate Earth's co-evolving

15    geosphere and biosphere. We therefore advocate a strategic, data-driven program for

16    accelerated scientific discovery.

17

18

19

22    _____

23    * E-mail address of corresponding author: rhazen@ciw.edu

24

2

## INTRODUCTION

25

26     Discovery in Earth sciences relies to a significant degree on induction and

27     deduction—classic approaches to reasoning that focus on the observation, modeling, and

28     ultimately (one hopes) predictive explanations of known patterns and phenomena in

29     nature. These powerful methods have proven successful in documenting and

30     comprehending many aspects of the natural world, but they are inherently inefficient at

31     discovering new complex patterns that require multivariate analysis of large datasets or

32     synthesis of diverse types of data. Consequently, recognition of such gradual global

33     processes as biological evolution by natural selection (Darwin 1859; Beddall 1968),

34     continental evolution by plate tectonics (Wood 1985; Hazen 2012), atmospheric and

35     ocean oxygenation by photosynthesis (Holland 1984; Canfield 2014), and climate change

36     (World Meteorological Organization 1989; Weart 2008) required decades of integrated

37     data synthesis preceding discovery and acceptance of critical Earth phenomena.

38     Today, the Earth and life sciences benefit from vast and ever-expanding data resources

39     in numerous disciplines—data that for the most part serve the needs of focused

40     communities of researchers. However, the potential now exists for a revolutionary

41     integration and synthesis of these diverse data resources, leading to an alternative

42     "abductive" approach to investigate Earth's co-evolving geosphere and biosphere. A

43     growing number of Earth scientists thus advocate a systematic, data-driven quest for the

44     accelerated discovery of hidden patterns in data resources from varied, interconnected

45     disciplines (Fayyad et al. 1996; Hazen et al. 2011; Keller and Schoene 2012; National

46     Science Foundation 2012; Bolukbasi et al. 2013; earthcube.org). Today's scientific

47     enterprises generate terabytes per day of new data, yet these vast resources are woefully

48  underutilized because they are not linked into a single platform (see, however, National

49  Science Foundation 2012). This "Outlooks" contribution examines strategies for linking

50  existing and new data resources, as well as methods for coupling integrated data

51  platforms with existing statistical analysis and visualization capabilities. Thus, we

52  envision a new kind of open-access "scientific instrument" that could transform the Earth

53  sciences.

54

55              **Deduction, Induction, and Abduction**

56      Most philosophers of science recognize two complementary modes of logical

57  reasoning by which many discoveries arise. By most definitions, deductive reasoning

58  begins with a general premise that is asserted to be true, and then draws specific

59  inferences from that generalization that must also be true. Thus:

60      • Earth's atmospheric oxygenation influenced the partitioning of redox-

61          sensitive elements.

62      • Molybdenum, rhenium, nickel, and cobalt are redox-sensitive elements.

63      • Therefore, we conclude by deduction that atmospheric oxygenation must

64          have influenced the partitioning of Mo. Re, Ni, and Co.

65  In deduction, specific conclusions represent a subset of the initial general premise.

66  Studies of the partitioning of redox-sensitive elements are thus conducted in the context

67  of well-established physical and chemical principles, and are not expected to yield

68  surprising or anomalous results that contradict the original premise. Such efforts are

69  critical to providing a solid foundation for scientific progress by filling in gaps in what

70  we know we don't know, but they do not usually represent the most efficient path to

71     discovering fundamentally new phenomena.

72       The complementary inductive mode of reasoning begins with observations of

73     particular instances of a generalization, which then lead to predictions of further instances

74     of the generalization (or to the generalization, itself). Thus:

75       • Each of the last 5 supercontinent cycles led to episodes of enhanced

76         mineralization during intervals of continental convergence.

77       • B, Be, Hg, and Mo are mineral-forming elements.

78       • Therefore, we predict by induction that B, Be, Hg, and Mo minerals will

79         display enhanced mineralization during intervals of continent convergence.

80     Unlike deduction, the specific predictions of induction are not necessarily contained

81     within the initial premise and thus they cannot follow with certainty. Because one starts

82     with instances of a generalization, and not an established premise, opportunities for

83     discovering unexpected or anomalous patterns may be enhanced. Thus, for example,

84     Hazen et al. (2012) found an anomalous absence of Hg mineralization during the

85     assembly of the Mesoproterozoic Rodinian supercontinent—an anomaly that parallels

86     emerging data from other studies (e.g., Huston et al. 2010; see below). The tradition in

87     Earth sciences (and crime novels) of collecting data to discriminate amongst multiple

88     working hypotheses (Chamberlin 1890) is inherently inductive in nature, and remains a

89     powerful strategy for discovery.

90       Most mineralogical research is firmly grounded in deduction and/or induction. Most

91     investigators, most of the time, start with an established deductive premise or an

92     inductive generalization consistent with observations about known phenomenon and then

93     collect new data to test the validity of one or more explanatory hypotheses, or to develop

94     new hypotheses.

95     These deductive and inductive efforts stand in contrast to "abduction", which is a form

96     of logical inference that begins with the accumulation of reliable data independently of a

97     premise or generalization. Analysis of these data, including statistical "data mining"

98     approaches, then point to previously unrecognized patterns and correlations, and

99     ultimately to the development of potentially new hypotheses to explain those patterns.

100    Discoveries that lead to "paradigm shifts"—for example, James Hutton's recognition of

101    gradual geological change and deep time (Hutton 1795), Charles Darwin's elucidation of

102    evolution by natural selection (Darwin 1859), and the collective development of the

103    concept of plate tectonics (Wood 1985)—tend to be intrinsically abductive in character,

104    even if the initial collection of data was motivated in a deductive/inductive context. Each

105    of these transformative discoveries required synthesis and integration of vast amounts of

106    diverse data resources accumulated over decades to articulate a new framing of the

107    natural world. Abduction thus provides a pathway to discovering what "we don't know

108    we don't know."

109

110                              **Data-Driven Discovery**

111    For most of the history of science abduction has proven a difficult and time-

112    consuming path to discovery. A lifetime of meticulous data collection and thoughtful

113    synthesis, at times amplified by creative intuition or blind luck, may be required to

114    recognize previously hidden patterns in diverse data. Only through decades of intimacy

115    with observations, and recognition of subtle quirks and idiosyncrasies in data, will some

116    significant patterns emerge from the noise. Such abductive discoveries do not easily

117   come to the impatient or distracted researcher, which should serve as an important

118   justification for the support of dedicated specialists who devote their lifetimes to a

119   focused scientific pursuit.

120      The development of large and expanding data resources, coupled with powerful

121   computation methods, has the potential to change the nature of abductive scientific

122   discovery. Advances in cyberinfrastructure are poised to integrate data from numerous

123   sources into semantically cohesive data platforms (Berner-Lee et al. 2001; Hey and

124   Trefethen 2005; Fox and Hendler 2009; Hey et al. 2009; McGuinness et al. 2009; Hazen

125   et al. 2011; Narock and Fox 2011). Furthermore, new and widely available statistical

126   methods and visualization procedures are providing the means to interrogate these data

127   resources in new ways and thus to tease out subtle correlations that are otherwise

128   inherently invisible to the human brain (Card et al. 1999; Hammer et al. 2001; Peter and

129   Shneiderman 2008; Fox and Hendler 2011; Kim et al. 2013).

130      Such data mining and discovery efforts that exploit large databases and enhanced data

131   interrogation techniques to seek patterns are much in the news, particularly with respect

132   to investment (Kovalerchuk and Vityaev 2000) and national security (Gellman and

133   Poitras 2013) applications. In science and medicine new data resources also have the

134   potential to reveal previously unrecognized phenomena. For example, seismological data

135   (from nuclear test ban verification efforts), coupled with ocean floor topography,

136   geochronology of ocean basalts, and paleomagnetism data, were critical in the discovery

137   of patterns that elucidated mechanisms of plate tectonics (e.g., Wood 1985). Today,

138   analyses of hidden patterns in genome databases to map viral evolution (Holmes 2007;

139   Lam et al. 2010) and statistical exploration of medical records to find potential causal

140  factors in pervasive diseases (Clos 2001; Berka et al. 2009) represent growing

141  applications of abductive strategies.

142  Similar opportunities await the mineralogist and petrologist. Earth materials scientists

143  have accumulated vast amounts of data on rocks, minerals, and geofluids, including their

144  major element, minor element, and isotopic compositions; optical, electrical, magnetic,

145  elastic, and other physical properties; atomic structures, as well as the variations of those

146  structures with pressure, temperature, and composition; petrologic context and associated

147  minerals; their ages; thermochemical parameters and phase relations; tectonic settings

148  and geologic context; and even their mineral-hosted microbial ecosystems (Table 1).

149  These data are complemented by resources on the evolution of paleoatmospheres and

150  paleooceans, geomicrobiology, paleontology, genomics and proteomics, paleotectonics,

151  paleomagnetism, and observations of other terrestrial planets and moons. The ultimate

152  goal of data-driven discovery is to create a single interoperable platform that offers

153  access to multiple, heterogeneous dimensions in a new "mineral data space." We can

154  envision a time when integrated data resources provide the key for discovering and

155  understanding numerous complementary aspects of Earth's evolution in space and time.

156  In spite of the promise of data-driven discovery, pitfalls abound. Data resources must

157  be approached with a firm grounding in chemical and physical principles; an awareness

158  of the meaning, quality, and sources of the data employed; and a keen sense of intuition.

159  Synthesis of unreliable or biased data from varied sources may lead to false or misleading

160  trends. For example, geochemical data employing different analytical instruments or

161  standardization procedures may display subtle systematic differences (Pyle et al. 2002;

162  Donovan et al. 2003). Other sources of bias reflect logistical factors: Recent studies of

163    mineral distributions in space and time (e.g., Hazen et al. 2012; Grew and Hazen 2014)

164    are invariably biased by the proximity of the most scrutinized deposits to major academic

165    institutions. Therefore, any interrogation of integrated data resources must be undertaken

166    within a framework of established deductive and inductive discovery.

167

168                                **Brute-Force Use Cases**

169        Integrated data resources for abductive discovery in mineralogy do not yet exist.

170    However, based on recent "brute-force use cases," we can be confident that previously

171    unrecognized patterns and correlations will emerge from the thoughtful integration and

172    evaluation of reliable data.

173        Brute-force use cases involve time-consuming, manual accumulation of relevant data,

174    either through literature searches or acquisition of new measurements. Such efforts have

175    been undertaken in many facets of the Earth sciences. Consider two recent examples from

176    the field of "mineral evolution"—studies of variations in the diversity and distribution of

177    the minerals of beryllium and mercury through deep time, which demonstrate the

178    potential of this concept as a means to recognize tectonic patterns; search for critical

179    resources; generate insights regarding the evolution of ocean and atmospheric chemistry;

180    and document subtle ongoing feedbacks among terrestrial life, weathering, soils, and

181    climate.

182        Grew and Hazen (2014), for example, accumulated age information for 122 Be

183    mineral localities, including the earliest known occurrences of all but 2 of the 112 known

184    minerals in which Be is an essential element. They collated these data from examination

185    of 300 references in 10 languages. This large dataset, assembled with minimal

186    preconceptions about what trends might emerge, revealed 5 significant episodes of Be

187    mineralization at approximately at 2600-2700, 1850-1750, 950-1050, 550-600, and 300

188    Ma—times in part associated with intervals of supercontinent assembly. Similar trends

189    are now emerging from data-intensive studies of granite pegmatites (Tkachev 2011), as

190    well as the ages of many thousands of individual detrital zircon crystals—large datasets

191    that add robustness to the interpretation of episodic mineralization over at least the past 3

192    billion years (Valley et al. 2005; Campbell and Allen 2008; Rino et al. 2008;

193    Hawksworth et al. 2010; Condie and Aster 2010; Condie et al. 2011; Voice et al. 2011).

194    Additional details in the Be dataset reveal other intriguing pulses of Be mineralization,

195    for example at ~1.3 Ga associated with extensional environments—a time not well

196    represented in the episodic zircon record.

197    In a similar effort, Hazen et al. (2012) surveyed 128 mercury mineral localities,

198    including the earliest known occurrences for 89 of the 90 known Hg species—a study

199    that required examination and synthesis of data in more than 400 references in a dozen

200    languages. Once again, this brute-force effort led to the discovery of previously hidden

201    patterns. In particular, 3 unexpected results emerged:

202    (1) The ages of almost all Hg mineral localities correlate with 4 episodes of

203        supercontinent assembly. Data fit to Gaussian curves at 2.69 ± 0.04, 1.81 ± 0.05,

204        0.53 ± 0.05, and 0.32 ± 0.07 Ga correlate with the assemblies of Kenorland, Nuna,

205        Pannotia, and Pangaea, respectively—patterns consistent with those observed for

206        zircon, molybdenite, and the minerals of Be and B.

207    (2) An as yet unexplained billion-year gap in Hg mineralization occurred between 1.8

208        and 0.8 billion years, an interval that included assembly of the supercontinent of

10

209    Rodinia. This interval may correlate with the innovation of microbial Hg

210    methylation, perhaps coupled with changes in ocean chemistry (Canfield 1998).

211    Alternatively, these data on Hg minerals may contribute to evidence that the

212    tectonic setting of Rodinian assembly differed from that of other supercontinents

213    (Cawood et al. 2009; Huston et al. 2010).

214     (3) The largest known Hg deposits from ~0.3 Ga are coeval with Carboniferous coal

215    measures, suggesting co-burial of organic carbon and Hg, with subsequent

216    hydrothermal mobilization and re-deposition. Thus, the mineralization of Hg—a

217    rare element with no biological function—is now coupled to the evolving terrestrial

218    biosphere.

219    This study was based entirely on published and web resources, yet it consumed more than

220    1 person-year of effort, mostly devoted to locating and evaluating previously published

221    data in sometimes obscure sources, as well as integrating those data with lists of minerals

222    approved by the International Mineralogical Association (Downs, 2006; rruff.info/ima)

223    and Hg locality information (principally in mindat.org).

224      These abductive discoveries, though focused on single rare elements with relatively

225    few localities, demonstrate the untapped potential for a new strategy of discovery based

226    on development and mining of enhanced data resources. Such brute-force studies of

227    Earth's near-surface rocks and minerals are by no means limited to the temporal

228    occurrence and aerial distribution of mineral species. For example, Faquahar et al. (2000,

229    2001, 2007) collated extensive data on sulfur isotopic fractionation in varied lithologies

230    versus age. They found remarkable mass-independent effects, a presumed consequence of

231    upper atmosphere photolysis of sulfur compounds, that are largely constrained to

232    formations > 2.25 Ga. They ascribed this finding to enhanced ozone shielding—a

233    conclusion heralded as the "smoking gun" for the timing of the Great Oxidation Event

234    and its irreversible transformation of atmospheric chemistry (Canfield 2014).

235    Large-scale community data resources, such as EarthChem/PetDB (Lehnert et al.

236    2000, 2007), are especially relevant for discoveries in statistical petrology, geochemistry,

237    and mineralogy (see http://www.earthchem.org/citations/petdb for examples). In one such

238    effort, Keller and Schoene (2012) employed a database of 70,000 analyses of continental

239    igneous rocks to discover evidence for significant lithospheric disruption at ~2.5 Ga—a

240    time just prior to the Great Oxidation Event. Their comprehensive overview of secular

241    variations in major and incompatible elements in basalt reveals a significant decrease in

242    mantle melt fraction at that time—a trend not obvious without a large and relatively

243    unbiased data set. Keller and Schoene (2012) concluded that atmospheric oxidation may

244    be linked in part to redox changes associated with crustal evolution. The availability of

245    large-scale, community supported, persistent, and quality-controlled data resources is

246    critical to the success of such endeavors.

247    Accumulations of mineral data also point to Earth's gradual subsurface oxidation.

248    Golden et al. (2013) gathered new and published trace element analyses of the rhenium

249    content of 422 molybdenite ($MoS_2$) specimens from 135 localities with known ages from

250    2.91 billion years to 6.3 million years. Rhenium is a redox-sensitive element that is

251    mobilized in its $Re^{7+}$ form only under relatively oxidized subsurface conditions. This

252    brute-force data effort revealed two statistically significant trends: (1) Systematic

253    increases in average and maximum trace concentrations of Re in molybdenite since 3.0

254    Ga point to enhanced oxidative weathering by subsurface fluids, and (2) episodic

255    molybdenum mineralization correlates with five intervals of supercontinent assembly

256    from ~2.7 Ga (Kenorland) to 300 Ma (Pangaea).

257        These and other examples demonstrate that brute-force methods have the potential to

258    reveal hidden correlations; numerous other trends in published mineralogical data are

259    undoubtedly awaiting discovery. However, brute-force data recovery methods are

260    inherently time-consuming and correspondingly inefficient. A far better strategy is to

261    develop and further enhance community supported data resources.

262

263              **A Strategy for Data-Driven Discovery in Mineralogy**

264        A strategy for data-driven discovery in mineralogy is emerging (Hazen et al. 2013).

265    The first steps relate to sustainable maintenance, expansion, and quality control of

266    existing databases for rocks, minerals, and other Earth materials (Table 1). The most

267    basic need and responsibility of any natural history discipline is the accurate, timely,

268    comprehensive, and accessible archiving of data on species. No task is more fundamental

269    to the long-term stability and integrity of a field, nor should such data management be left

270    to the unfunded good will of individuals, no matter how skilled and well intentioned they

271    may be. Some data resources such as EarthChem and Volcanoes (Table 1) enjoy

272    significant, if not guaranteed long-term, Federal support. However, it is astonishing that

273    the official web-based list of approved mineral species (rruff.info/ima), which is freely

274    available and widely used by the international community, has no long-term institutional

275    home or financial support. And, until recently, the server for LEPR and MELTS—widely

276    used open-access resources for thermochemistry—resided in the bedroom of its founder,

13

277    Mark Ghiorso. The international Earth materials community, therefore, needs to initiate

278    action on several fronts:

279      (1) Encourage professional journals to adopt policies and collaborations by which all

280          newly published data will be deposited in approved, sustainable, quality-controlled,

281          open-access sources.

282      (2) Endorse the International Geo Sample Number (IGSN) procedures of the System

283          for Earth Sample Registration (Lehnert and Klump 2008; www.geosamples.org)

284          and require use of this system before publication.

285      (3) Continue to encourage and support efforts to transfer previously published data

286          into open-access repositories.

287      (4) Identify and encourage recovery of "dark data" resources, including unpublished

288          hard copy and electronic format data accumulated by individuals. Encourage

289          publication of these data resources through electronic data journals with digital

290          object identifier (doi) information.

291      (5) Create active and engaged user communities to ensure quality control of data

292          resources, which must be properly vetted prior to incorporation into open-access

293          sources.

294      (6) Establish data publication procedures and data citation policies that ensure proper

295          credit and motivation for data producers.

296      (7) Identify and exploit sources of funding for work being done now and for long-term

297          institutional support of critical data resources.

298    Note that several of these steps will require both institutional and cultural changes within

299    the Earth sciences community. Effective change will take time and effort, but we can

300  anticipate a time when the larger Earth materials community recognizes the critical

301  importance of shared, high-quality data resources and underscores the responsibility of all

302  researchers to contribute to this infrastructure.

303  The second facet of fostering data-driven discovery in Earth materials research

304  involves integrating existing databases into a larger Earth Materials Data Infrastructure

305  (Hazen et al. 2013). Ultimately, we can envision linking separately maintained data

306  resources into a federated, centrally-governed data framework in which diverse data

307  resources are semantically compatible and linked (see, for example, National Science

308  Foundation 2012; earthcube.org). Eventually, a more expansive opportunity lies in the

309  integration of Earth materials data with other complementary disciplines, including

310  paleobiology, proteomics, paleotectonics, and planetary sciences.

311  To accomplish this vision, we need to identify and integrate key data resources, while

312  providing computational tools that can be used to select, analyze, and visualize data.

313  Ultimately, a comprehensive Earth materials data infrastructure could be linked to

314  artificial intelligence and machine learning capabilities to accelerate data-driven

315  discovery. Cyberinfrastructure programs such as EarthCube, which enjoy significant

316  community support as well as Federal funding, are moving scientific research in these

317  bold new directions. However, the mineralogy-petrology research community needs to

318  increase its commitment to these efforts if we are to take maximum advantage of

319  emerging opportunities.

320

321                                        **Conclusions**

322     Two significant and novel impacts are likely to result from a program of abductive

323     discovery in mineralogy. First, unanticipated mineralogical discoveries will be made.

324     New phenomena related to mineral crystal chemistry, trace element and isotope

325     distributions, mineral associations and geologic context, fluid-rock interactions, and

326     interactions with the biosphere, all in the framework of geological space and time, are

327     certain to emerge. The abductive approach, coupled with more traditional hypothesis-

328     driven inquiries, will inevitably lead to discovery of new patterns in nature.

329     Second, and more significant for the future of mineralogy and petrology, data-driven

330     discovery provides a model for $21^{st}$-century science that explicitly recognizes the power

331     of abduction and exploits opportunities represented by the explosion of Earth science

332     data. Such a strategy in no way subsumes deduction and induction; indeed, the abductive

333     approach explicitly relies upon the accumulation of traditional measurements and

334     observations. Abduction amplifies the vast amounts of data inspired by deductive and

335     inductive discovery. Ultimately, when data from petrology, mineralogy, geochemistry,

336     paleontology,     geodynamics,     proteomics,     irreversible     thermodynamics,     and

337     geochronology are integrated with newly adapted statistical analysis and visualization

338     capabilities, we will enjoy a wholly new kind of "scientific instrument"—an open-access

339     engine of discovery that could transform the Earth sciences.

340

16

353

354                                        **References**

355

356   Alroy, J. (2010) The shifting balance of diversity among major animal groups. Science,

357        329, 1191-1194.

358   Alroy, J., Aberhan, M., Bottjer, D.J., Foote, M., Fürsich, F.T., Harries, P.J., Hendy,

359        A.J.W., Holland, S.M., Ivany, L.C., Kiessling, W., Kosnik, M.A., Marshall, C.R.,

360        McGowan, A.J., Miller, A.I., Olszewski, T.D., Patzkowsky, M.E., Peters, S.E., Villier,

361        L., Wagner, P.J., Bonuso, N., Borkow, P.S., Brenneis, B., Clapham, M.E., Fall, L.M.,

362        Ferguson, C.A., Hanson, V.L., Krug, A.Z., Layou, K.M. Leckey, E.H., Nürnberg, S.,

363        Powers, C.M., Sessa, J.A., Simpson, C., Tomasovych, A., and Visaggi, C.C. (2008)

364        Phanerozoic trends in the global diversity of marine invertebrates. Science, 321, 97-

365        100.

366   Beddall, B. G. (1968). Wallace, Darwin, and the theory of natural selection. Journal of

367        the History of Biology, **1,** 261–323.

368   Berka, P., Rauch, J., and Djamel, A.Z. [Editors] (2009) Data Mining and Medical

369        Knowledge Management: Cases and Applications. Hershey, Pennsylvania: IGI Global.

370   Berner-Lee, T., Hendler, J., and Lassila, O. (2001) The semantic web. Scientific

371        American, 284(5), 34-43.

372   Bolukbasi, B., Berente, N., Cutcher-Gershenfeld, J., Dechurch, L., Flint, C., Haberman,

373        M., King, J.L., Knight, E., Lawrence, B., Masella, E., McElroy, C., Mittleman, B.,

374        Nolan, M., Radik, M., Shin, N., Thompson, C.A., Winter, S., Zaslavsky, Allison,

375        M.L., Arctur, D., Arrigo, J., Aufdenkampe, A.K., Bass, J., Crowell, J., Daniels, M.,

376        Diggs, S., Duffy, C., Gil, Y., Gomez, B., Graves, S., Hazen, R., Hsu, L., Kinkade, D.,

377    Lehnert, K., Marone, C., Middleton, D., Noren, A., Paerthree, G., Ramamurthy, M.,

378    Robinson, E., Percivall, G., Richard, S., Suarez, C., and Walker, D. (2013) Open data:

379    Crediting a culture of cooperation. Science, 342, 1041-4042.

380    DOI:10.1126/science.342.6162.1041-b

381 Campbell, I.H., and Allen, C.M. (2008) Formation of supercontinents linked to increases

382    in atmospheric oxygen. Nature Geoscience, 1, 554-558.

383 Canfield, D.E. (1998) A new model for Proterozoic ocean chemistry. Nature, 396, 450-

384    453.

385 Canfield, D.E. (2014) Oxygen: A Four Billion Year History. Princeton, NJ: Princeton

386    University Press.

387 Card, S.K., Mackinlay, J.D., and Shneiderman, B. (1999) Reading in Information

388    Visualization: Using Vision to Think. San Francisco, California: Morgan Kaufmann.

389 Cawood, P.A., Kroner, A., Collins, W.J., Kusky, T.W., Mooney, W.D., and Windley,

390    B.F. (2009) Accretionary orogens through Earth history. Geological Society [London]

391    Special Publication, 318, 1−36.

392 Chamberlin, T.C. (1890) The method of multiple working hypotheses. Science, 15, 92-

393    96.

394 Clos, K.J. [Editor] (2001) Medical Data Mining and Knowledge Discovery. Dortrecht,

395    Netherlands: Springer.

396 Condie, K.C., and Aster, R.C. (2010) Episodic zircon age spectra of orogenic granitoids:

397    The supercontinent connection and continental growth. Precambrian Research, 180,

398    227-236.

399 Condie, K.C., Bickford, M.E., Aster, R.C., Belousova, E., and Scholl, D.W. (2011)

19

400    Episodic zircon ages, Hf isotopic composition, and the preservation rate of continental

401    crust. Geological Society of America Bulletin, 123, 951-957.

402  Darwin, C. (1859) On the Origin of Species. London: John Murray.

403  Donovan, J.J., Hanchar, J.M., Picollo, P.M., Schrier, M.D., Boatner, L.A., and

404    Jarosewich, E. (2003) A Re-examination of the rare-earth-element orthophosphate

405    standards in use for elecrtron-microprobe analysis. Canadian Mineralogist. 41, 221-

406    232.

407  Downs, R.T. (2006) The RRUFF Project: an integrated study of the chemistry,

408    crystallography, Raman and infrared spectroscopy of minerals. Program and Abstracts

409    of the 19th General Meeting of the International Mineralogical Association in Kobe,

410    Japan. O03-13.

411  Downs, R.T., and Hall-Wallace, M. (2003) The American Mineralogist Crystal Structure

412    Database. American Mineralogist, 88, 247-250.

413  Fox, P., and Hendler, J. (2009) Semantic eScience: Encoding meaning in next-generation

414    digitally enhanced science. In Hey, T., Tansley, S., Tolle, K. [Editors] The Fourth

415    Paradigm: Data-Intensive Scientific Discovery. Redland, WA: Microsoft External

416    Research, pp. 145-150.

417  Farquhar, J., Bao, H., and Thiemens, M.H. (2000) Atmospheric Influence of Earth's

418    Earliest Sulfur Cycle. Science, 289, 756-758.

419  Farquhar, J., Savarino, I., Airieau, S., and Thiemens, M.H. (2001) Observations of

420    wavelength- sensitive, mass-independent sulfur isotope effects during $SO_2$ photolysis:

421    Implications for the early atmosphere. Journal of Geophysical Research, 106, 1-11.

422    Farquhar, J., Peters, M., Johnston, D.T., Strauss, H., Masterson, A., Wiechert, U., and

423         Kaufman, A.J. (2007) Isotopic evidence for mesoarchean anoxia and changing

424         atmospheric sulphur chemistry. Nature, 449, 706-709.

425    Fayyad, U., Piatetsky-Shapiro, G., and Smyth, P. (1996) From data mining to knowledge

426         discovery in databases. AI Magazine, Fall 1996, 37-54.

427    Fox, P., and Hendler, J. (2011) Changing the equation on scientific data visualization.

428         Science, 331, 705-708.

429    Gellman, B. and Poitras, L. (2013). U. S. intelligence mining data from nine U.S. Internet

430         companies in broad secret program. The Washington Post, June 7, 2013, A1.

431    Ghiorso, M.S., and Sack, R.O. (1995) Chemical Mass Transfer in Magmatic Processes

432         IV. A revised and internally consistent thermodynamic model for the interpolation and

433         extrapolation of liquid-solid equilibria in magmatic systems at elevated temperatures

434         and pressures. Contributions to Mineralogy and Petrology, 119, 197-212.

435    Ghiorso, M.S., Hirschmann, M.M., Reiners, P.W., and Kress, V.C. III (2002) The

436         pMELTS: A revision of MELTS for improved calculation of phase relations and

437         major element partitioning related to partial melting of the mantle to 3

438         GPa. Geochemistry Geophysics Geosystems, 3, 10.1029/2001GC000217.

439    Grazulis, S., Daskevic, A., Merkys, A., Chateigner, D., Lutterotti, L., Quiros, M.,

440         Serebryanaya, N.R., Moeck, P., Downs, R.T., and Le Bail, A. (2012) Crystallography

441         Open Database (COD): an open-access collection of crystal structures and platforms

442         for      world-wide      collaboration. Nucleic      Acids      Research, 40,      D420-D427.

443         doi:10.1093/nar/gkr900

444   Grew, E.S., and Hazen, R.M. (2014) Beryllium mineral evolution. American

445        Mineralogist, in press.

446   Hammer, Ø., Harper, D.A.T., and Ryan, P.D. (2001) PAST: Paleontological statistical

447        software package for education and data analysis. Paleo-Electronica.org, issue 1.

448   Hawkesworth, C.J., Dhuime, B., Pietranik, A.B., Kemp, A.I.S., and Storey, C.D. (2010)

449        The generation and evolution of continental crust: Journal of the Geological Society,

450        167, 229-248.

451   Hazen, R.M. (2012) The Story of Earth. New York: Viking-Penguin.

452   Hazen R.M., Bekker A., Bish D.L., Bleeker W., Downs R.T., Farquhar J., Ferry J.M.,

453        Grew E.S., Knoll A.H., Papineau D., Ralph J.P., Sverjensky D.A., and Valley J.W.

454        (2011) Needs and opportunities in mineral evolution research. American Mineralogist,

455        96, 953-963.

456   Hazen, R.M., Golden, J., Downs, R.T., Hystad, G., Grew, E.S., Azzolini, D, and

457        Sverjensky, D.A. (2012) Mercury (Hg) mineral evolution: A mineralogical record of

458        supercontinent assembly, changing ocean geochemistry, and the emerging terrestrial

459        biosphere. American Mineralogist, 97, 1013-1042.

460   Hazen, R.M., Cotrell, E., Downs, R.T., Fox, P., Ghiorso, M., Lehnert, K., Saxena, S., and

461        Spear, F. (2013) Report of the Chair of the Ad Hoc Committee on Earth Materials

462        Data, To the First 2013 Mineralogical Society of America Council Meeting, 23 April

463        2013, 6 pp.

464   Hey, T., and Trefethen, A.E. (2005) Cyberinfrastructure for e-Science. Science, 308, 817-

465        821.

466  Hey, T., Tansley, S., and Tolle, K. [Editors] (2009) The Fourth Paradigm: Data-Intensive

467      Scientific Discovery. Redland, WA: Microsoft External Research.

468  Holland, H.D. (1984) The Chemical Evolution of the Atmosphere and Oceans. Princeton,

469      NJ: Princeton University Press.

470  Holland, T.J.B., and Powell, R. (1998) An internally consistent thermodynamic data set

471      for phases of petrological interest. Journal of Metamorphic Petrology, 16, 309–343.

472  Holmes, E.C. (2007) Viral evolution in the genomic age. PLOS Biology, DOI:

473      10.1371/journal.pbio.0050278

474  Huston, D.L., Pehrsson, S., Eglington, B.M., and Zaw, K. (2010) The geology and

475      metallogeny of volcanic-hosted massive sulfide deposits: Variations through geologic

476      time and with tectonic setting. Economic Geology, 106, 571-591.

477  James Hutton (1795) Theory of the Earth; with Proofs and Illustrations. Edinburgh:

478      Creech. 2 volumes.

479  Keller, B., and Schoene (2012) Statistical geochemistry reveals disruption in secular

480      lithospheric evolution about 2.5 Gya ago. Nature, 485, 490-493.

481  Kim, J.D., Senn, S., Harel, A., Jelen, B.I., and Falkowski, P.G. (2013) Discovering the

482      electronic circuit diagram of life: structural relationships among transition metal

483      binding sites in oxidoreductases. Philosophical Transactions of the Royal Society, B

484      368, 20120257.

485  Kovalerchuk, B. and Vityaev, E. (2000) Data Mining in Finace: Advances in Relational

486      and Hybrid Methods. New York: Kluwer Academic Publishers.

487  Lam, T.T., Hon, C.C., and Tang, J.W. (2010) Use of phylogenetics in the molecular

488      epidemiology and evolutionary studies of viral infections. Critical Reviews in Clinical

489      Laboratory Sciences, 47, 5–49.

490  Lehnert, K.A., and Klump, J. (2008) Facilitating research in mantle petrology with

491      geoinformatics. 9th International Kimberlite Conference Extended Abstracts, 91KC-A-

492      00250.

493  Lehnert, K.A., Su, Y., Langmuir, C.H., Sarbas, B., and Nohl, U. (2000) A global

494      geochemical database structure for rocks. Geochemistry Geophysics Geosystems 1.

495  Lehnert, K.A., Walker, D., and Sarbas, B. (2007) EarthChem: A geochemistry data

496      network. Geochimica et Cosmochimica Acta, 71, A559.

497  McGuinness, D.L., Fox, P.A., Brodaric, B., and Kendall, E. (2009) The emerging field of

498      semantic scientific knowledge integration. IEEE Intelligent Systems, 24, 25-26.

499  Mutschler, F.E., Rougon, D.J., Lavin, O.P., and R.D. Hughes (1981) PETROS version

500      6.1 Worldwide Databank of Major Element Chemical Analyses of Igneous Rocks.

501      National Geophysical Data Center, NOAA. doi:10.7289/V5QN64NM.

502  Narock, T., and Fox, P.A. (2011) From science to e-science to semantic e-science: A

503      heliophysics case study. Computers & Geosciences. doi: 10.1016/j.cageo.2011.11.018

504  National Science Foundation (2012) A Community Roadmap for Earthcube Data:

505      Discovery, Access, and Mining. Arlington, Virginia: National Science Foundation, 38

506      p.

507  Perer, A., and Shneiderman, B. (2008) Integrating statistics and visualization: Case

508      studies of gaining clarity during exploratory data analysis. ACM Conference on

509      Human Factors in Computing Systems, Florence, Italy.

510    Pyle, J.M., Spear, F.S., and Wark, D.A. (2002) Electron microprobe analysis of REE in

511        apatite, monazite and xenotime: Protocols and pitfalls. Reviews in Mineralogy and

512        Geochemistry, 48, 337-362.

513    Rino, S., Kon, Y., Sato, W., Maruyama, S., Santosh, M., and Zhao, D. (2008) The

514        Grenvillian and Pan-African orogens: world's largest orogenies through geologic time,

515        and their implications on the origin of superplumes. Gondwana Research, 14, 51-72.

516    Sepkowski, D. (2012) Towards "A natural history of data": Evolving practices and

517        epistemologies of data in paleontology, 1800-2000. Journal of the History of Biology.

518        DOI 10.1007/s10739-012-9336-6

519    Stixrude, L., and Lithgow-Bertelloni, C. (2005) Thermodynamics of mantle minerals – I.

520        Physical properties. Geophysical Journal International, 162, 610–632.

521    Stixrude, L., and Lithgow-Bertelloni, C. (2011) Thermodynamics of mantle minerals – II.

522        Phase equilibria, Geophysical Journal International, 184, 1180–1213

523    Tkachev, A.V. (2011) Evolution of metallogeny of granitic pegmatites associated with

524        orogens throughout geological time. Geological Society of London, Special

525        Publications 350, 7-23.

526    Valley, J.W., Lackey, J.S., Cavosie, A.J., Clechenko, C.C., Spicuzza, M.J., Basei,

527        M.A.S., Bindeman, I.N., Ferreira, V.P., Sial, A.N., King, E.M., Peck, W.H., Sinha,

528        A.K., and Wei, C.S. (2005) 4.4 billion years of crustal maturation: oxygen isotope

529        ratios of magmatic zircon. Contributions to Mineralogy and Petrology, 150, 561-580.

530    Voice, P.J., Kowalewski, M., and Eriksson, K.A. (2011) Quantifying the timing and rate

531        of crustal evolution: Global compilation of radiometrically dated detrital zircon grains.

532        The Journal of Geology, 119, 109-126.

25

533    Weart, S.R. (2008) The Discovery of Global Warming: Revised and Expanded Edition.

534        Cambridge, Massachusetts: Harvard University Press.

535    Wood, R.M. (1985) The Dark Side of the Earth. London: George Allen and Unwin.

536    World Meteorological Organization (1989) The Changing Atmosphere: Implications for

537        Global Security, Toronto, Canada, 27-30 June 1988: Conference Proceedings. Geneva:

538        Secretariat of the World Meteorological Organization.

539

540    Table 1. Selected open-access data resources for Earth materials research.

541

| Web address | Content | ref |
|---|---|---|
| rruff.info | Mineral species and properties | 1 |
| rruff.info/ima | IMA official mineral list | |
| rruff.geo.arizona.edu/AMS/amcsd.php | Mineral crystal structures | 2 |
| smmp.net/IMA-CM/ctms.htm | IMA list of type minerals | |
| www.crystallography.net | Crystal structure data | 3 |
| http://cod.iutcaen.unicaen.fr | Powder diffraction data | |
| http://database.iem.ac.ru/mincryst/ | Mineral crystal structures | |
| mindat.org | Mineral localities/associations/properties | |
| webmineral.com | Mineral species/properties | |
| athena.unige.ch/athena/ | Mineral species/properties | |
| http://abulafia.mt.ic.ac.uk/shannon/ptable.php | Ionic radii table | |
| http://minerals.gps.caltech.edu/FILES/raman/Caltech_data/index.htm | | |
| | Spectroscopic data | |
| earthref.org | geochemistry/geomagnetism data | |
| geokem.com | igneous rock chemistry | |
| georoc.mpch-mainz.gwdg.de | rock geochemistry | |
| metpetdb.rpi.edu | metamorphic petrology | |
| navdat.org | igneous rocks of North America | |
| earthchem.org | geochemistry, geochronology, petrology | |
| earthchem.org/petdb | petrology | |

27

| | | |
|---|---|---|
| 563 | ngdc.noaa.gov/mgg/geology/petros.html | igneous rock geochemistry 4 |
| 564 | http://volcano.si.edu | volcanoes and eruptions |
| 565 | iza-structure.org/databases/ | zeolite crystal structures |
| 566 | melts.ofm-research.org | thermodynamic modeling 5,6 |
| 567 | lepr.ofm-research.org | experimental data |
| 568 | metamorph.geo.uni-mainz.de/thermocalc/ | thermochemical modeling 7 |
| 569 | phaseplot.com/Phase_Plot/Contents.html | mantle phase equilibria modeling 8 |
| 570 | vamps.mbl.edu/portals/deep_carbon/cdl.php | subsurface microbial ecosystems |

571 _____

572  1. Downs (2006); 2. Downs and Wallace (2003); 3. Grazulis et al. (2012); 4. Mutschler et

573  al. (1981); 5. Ghiorso and Sack (1995); 6. Ghiorso et al. (2002); 7. Holland and Powell

574  (1998); 8. Stixrude and Lithgow-Bertelloni (2005, 2011).

575 _____

576