# A skeleton data model for geochemical databases at the British Columbia Geological Survey

Tian Han[1, a], Alexei S. Rukhlov[1], Janet M. Riddell[1], and Travis Ferbey[1]

[1] British Columbia Geological Survey, Ministry of Energy, Mines and Petroleum Resources, Victoria, BC, V8W 9N3
[a] corresponding author: Tian.Han@gov.bc.ca

## Abstract

Data modelling is a key component to developing any database. In the past, the British Columbia Geological Survey has modelled geochemical data sets individually. Herein we propose a skeleton data model capable of capturing and representing the commonalities of individual sets by focusing on the entities, attributes, and relationships common to most geochemical data sets. We use this skeleton data model to update existing models for four province-wide data sets (lithogeochemical, regional drainage, till, and coal ash), capturing the unique characteristics of each. Applying the skeleton data model streamlines data handling steps, from data compilation to product generation, and establishes a reliable flow for managing geochemical data at the British Columbia Geological Survey.

**Keywords:** Geochemical data model, database, entities, attributes, relationships, data lifecycle activities, operation, query, SQL, data access, analytical methods, chemical element abundance
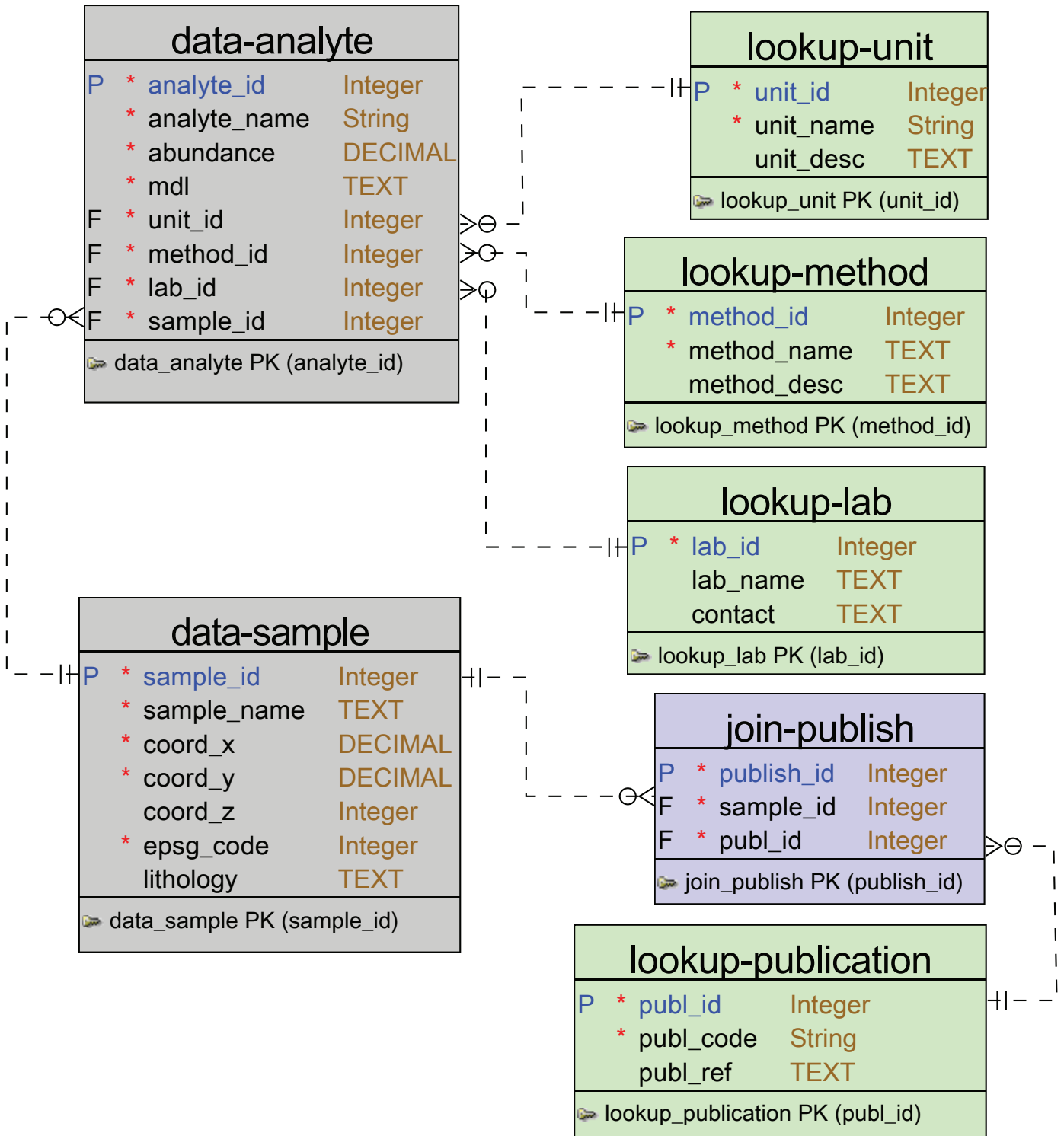
## 1. Introduction

The British Colombia Geological Survey (BCGS) is currently the custodian of four province-wide geochemical data sets: lithogeochemical; regional drainage geochemical surveys (RGS); till; and coal ash. The lithogeochemical data set contains analyses from bedrock samples; the RGS data set contains analyses from water and stream-, lake- and moss mat- sediment samples; the till geochemical data set contains analyses from subglacial till samples; and the coal ash data set contains analyses of the inorganic residue remaining after coal combustion. These data sets were extracted from reports produced or archived by the BCGS in PDF, Excel, or ASCII formats. The BCGS has begun systematically compiling and migrating these geochemical data sets to relational databases, where data are centrally maintained and managed for easy access, efficient update, enhanced consistency, effective quality control, and long-term storage. Previously, we created individual data models for each of the four province-wide data sets (Han et al., 2016, 2017; Bustard et al., 2017; Riddell and Han, 2017). However, although the data sets are derived from samples of different media, they all record values of abundances of chemical elements or compounds and thus share commonalities that can be expressed in terms of generic entities, common attributes, and intrinsic relationships. This led us to develop a skeleton data model to capture these commonalities across all geochemical data sets that could be used as a base for developing data set-specific models. These data set-specific models are created by customizing the skeleton model to incorporate the unique characteristics of the corresponding data sets, rather than by building an entirely new model. This approach helps standardize geochemical data attribution, management, reconciliation, publication, and quality control.

In this paper, we present the skeleton data model, describe how it is adapted for each of the four province-wide geochemical data sets, present the corresponding detailed data model, and outline the post-data modeling work, including the streamlined process built around these databases to support geochemical data lifecycle activities.

## 2. Developing a skeleton data model for geochemical data

Database development is typically done in four stages: requirement analysis; logical design; implementation; and database population (Connolly and Begg, 1999). The result of logical design is a data model conceived to represent the data of interest in a database environment. A data model is obtained through a process called data modeling by: 1) determining the data entities; 2) defining the attributes for each data entity; and 3) resolving the relationship between data entities. Each of the lithogeochemical, RGS, till, and coal ash geochemical data sets has multiple entities, attributes, and relationships. Some are unique to a specific data set, whereas others are common to all data sets. Below we consider the common ones and outline a skeleton data model capable of representing them all (Fig. 1). To complement the 'structural view' portrayed by the skeleton data model (Fig. 1), a 'data view' (Fig. 2) shows the appearance of the data model once it is populated with real data. To avoid confusion in the following discussion, we present entity names in boldface and attribute names in italics.

125

Geological Fieldwork 2018, British Columbia Ministry of Energy, Mines and Petroleum Resources, British Columbia Geological Survey Paper 2019-01

**Fig. 1.** Skeleton data model for geochemical data. Each of the seven boxes represents an entity. The dashed lines connecting the boxes depict entity relationships. The name of each entity is shown in boldface; attribute names are on the left, related data types on the right. The bottom row of each box is the attribute that functions as the primary key of the entity. Mandatory attributes are indicated by an asterisk.
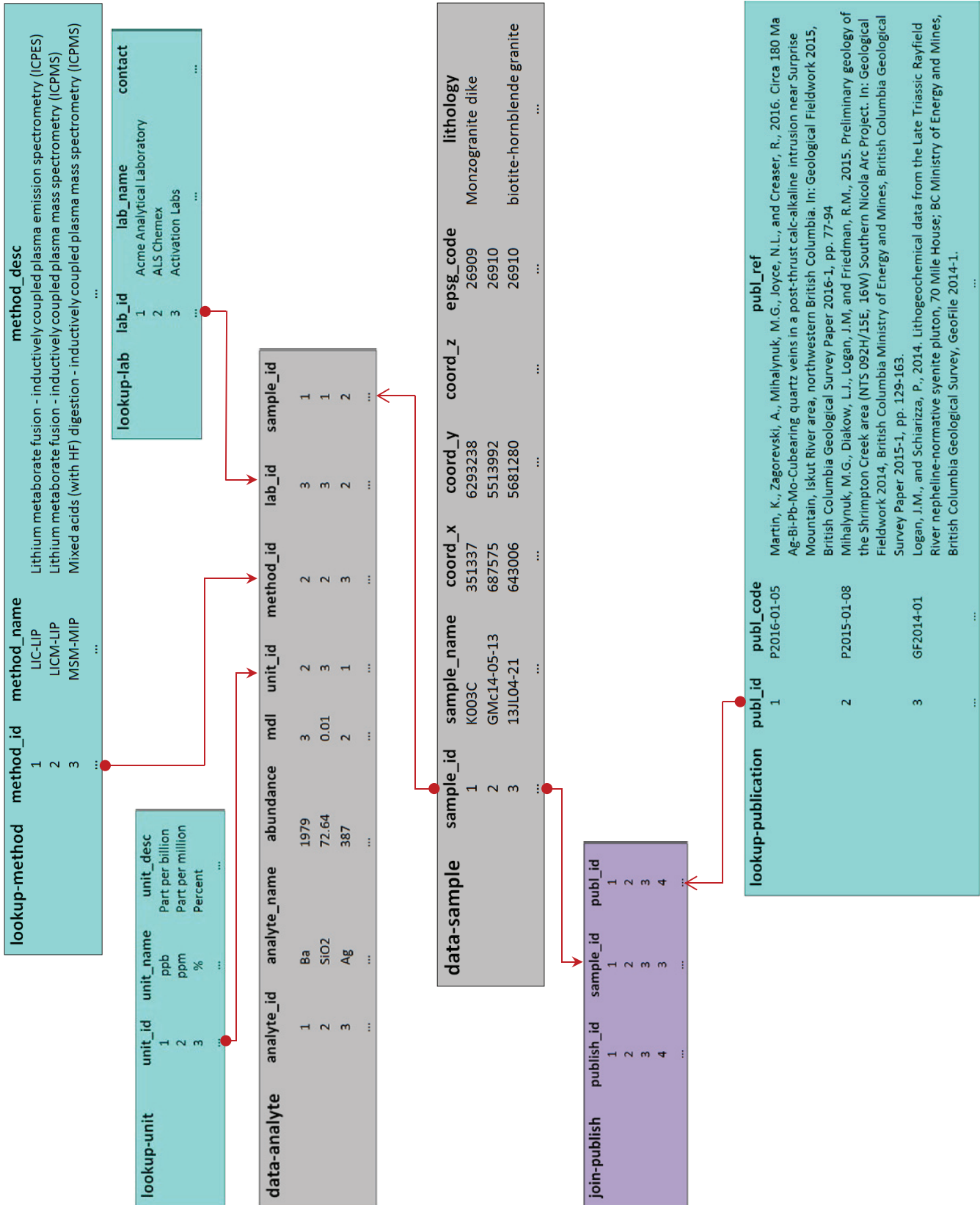
**Fig. 2.** Skeleton data model populated with real data in 'data view'.

127

Geological Fieldwork 2018, British Columbia Ministry of Energy, Mines and Petroleum Resources, British Columbia Geological Survey Paper 2019-01

## 2.1. Entities

In general, a geochemical data set comprises two generic entities: sample and analyte, and several related auxiliary entities. The sample entity comprises a suite of common attributes, including sample name, location, geology, lithology, and other properties. In the case where one field sample is split into several fractions, each split is a sub-sample, which is treated as a formal record in the sample entity and is assigned with a unique *sample_id*. These sub-samples can be further named in the *sample_name* attribute so that the 'parent-children' relationship between a field sample and its sub-samples is made obvious. The analyte entity captures analytical values representing concentrations of chemical elements or compounds. The common attributes of the analyte entity are analytical value, unit, analytical method, and the laboratory responsible for the analysis. These two entities and their associated attributes are typically common to all geochemical data sets. The skeleton data model (Fig.1) has the sample and analyte entities, along with five other auxiliary entities, where each entity is depicted as a table with a list of entries for its attributes and data types.

Among the five auxiliary entities included in the skeleton data model (Fig. 1), the **lookup-unit** holds all common units that can be assumed by the analyte abundance values; the **lookup-method** captures all the analytical methods used for sample analysis; the **lookup-lab** maintains the analytical laboratories and contact information; and the **lookup-publication** provides source references for samples included in the database. All these entities have names starting with 'lookup', indicating they function as lookup tables. Though these entities could be combined with either the **data-sample** or **data-analyte** entities, keeping them separate is preferred for categorizing logic data, normalizing between entities, and enhancing data integrity and performance. For example, if an analytical laboratory changes its name from X to Y and the laboratory information is stored in the **data-analyte** entity, we would have to update all the records of analytes that were done by X. Missing any one of them would compromise data integrity. By keeping laboratory information separately as shown in Figure 1, we only have to update one record in the **lookup-lab** entity.

## 2.2. Attributes

The sample entity, named **data-sample** in the skeleton data model, contains a list of nine attributes. These attributes are *sample_id* for record identification; *sample_name* for sample name; *coord_x* and *coord_y* for sample location coordinates; *coord_z* for sample elevation; *epsg_code* for geospatial reference system code defined by the European Petroleum Survey Group (EPSG; IOGP, 2018); *collect_date* for date of sample collection; and *lithology* and *sample_desc*. Using EPSG code enables us to represent sample coordinates concisely in their original spatial reference systems and avoid unnecessary re-projection calculations. The analyte entity, named **data-analyte** in the skeleton data model, has eight attributes, of which *analyte_id* is for record identification; *analyte_name* is

for analyte name; *abundance* is for analyte abundance value; and *mdl* is for minimum detection limit.

## 2.3. Relationships

Having determined entities and their associated attributes, we then resolve relationships between entities. For each entity of the skeleton data model, the first attribute is created for record identification with a name that ends with 'id'. Although without real physical meaning, this attribute is called the primary key, and is a unique identifier that can be used to locate specific records and facilitate efficient joining between entities. Other attributes with names ending with 'id' in each entity (Fig. 1) are used for linking related entities. These attributes are called foreign keys. A foreign key in one entity is commonly the primary key in another. For example, *sample_id* is a foreign key in **data-analyte** entity but the primary key in **data-sample** entity.

The relationship between the entities **data-sample** and **data-analyte** is determined by the *sample-id* in both. It is a 'one-to-many relationship', meaning that one sample may have many analytes. This is common in geochemical data sets because samples are routinely analyzed for multiple elements. One-to-many relationships also exist between **lookup-unit** and **data-analyte** defined by *unit_id*; **lookup-method** and **data-analyte** by *method_id*; and **lookup-lab** and **data-analyte** by *lab_id*. These relationships are depicted using a 'crow foot' symbol (Connolly and Begg, 1999) in Figure 1.

The **data-sample** and **lookup-publication** entities commonly have a 'many-to-many relationship', such as where one sample is reported in multiple publications or, as is generally the case, a single publication reports results from multiple samples. It would be difficult and inefficient to implement this relationship directly because it would result in many duplicated or blank attributes. As a result, a special entity, **join-publish**, is introduced and is embedded between **data-sample** and **lookup-publication**. It joins the two (hence the name starting with join-) and turns the many-to-many relationship into two one-to-many relationships; one between **data-sample** and **join-publish** and the other between **lookup-publication** and **join-publish**.

## 2.4. Query support

Our skeleton data model also considers user query requirements. The model and derived databases are designed to accommodate common query requirements such as looking up records of samples or analytes that satisfy certain criteria. For example, a user may want to see samples with Ag concentrations above 50 ppm that were analyzed with AAS (aqua regia-cold vapour atomic absorption spectrometry). Or the user may want to identify samples with Au and As concentrations (analyzed using Instrumental Neutron Activation) both above 95[th] percentile. The model is able to accommodate unique identification and extraction of samples efficiently with simple Structured Query Language (SQL) statements.

The model also supports spatial queries. Geochemical

128

Geological Fieldwork 2018, British Columbia Ministry of Energy, Mines and Petroleum Resources, British Columbia Geological Survey Paper 2019-01

sample sites are treated as points with coordinates specified in geospatial reference systems using EPSG codes (a standard in geospatial information technology and recognized by the Open Geospatial Consortium). Sample coordinates can be extracted and used for spatial visualization and further examination against other geological data sets.

## 3. Adapting the skeleton model to suit different geochemical data sets

The skeleton data model presented in Figure 1 aligns well with the Open Geoscience data model (Granitto et al, 2012; Watson and Evans, 2012). It is our intention to use it as the framework to model all geochemical data sets currently hosted by the BCGS. It is a base that can be built upon to address the unique entities, attributes, and relationships existing in the corresponding data sets. Below we present the data model for each of the four province-wide data sets and discuss the customization details.

### 3.1. Data model for lithogeochemical data set

The BCGS lithogeochemical data set consists of more than 11,000 samples with about a quarter million determinations generated by 26 different analytical methods carried out at 21 laboratories. The data were compiled from reports published since 1986 by BCGS geoscientists and research partners from the Geological Survey of Canada (GSC) and Canadian universities.

We reconstructed the data model used by Han et al (2016) based on the skeleton model. The updated lithogeochemical data model (Fig. 3) retains all components of the skeleton model with the addition of several lithogeochemistry-specific attributes and two entities. A **lookup-geologist** entity was added to identify who collected a sample and a **lookup-preparation** entity was included to record how samples were prepared before analysis.

### 3.2. Data model for RGS geochemical data set

The regional geochemical survey data set includes results from drainage samples (stream-, lake-, and moss mat-sediment, and water), collected by the Geological Survey of Canada, BCGS, and Geoscience BC since 1976. About 80% of the province has been sampled at a density of one sample per 7 to 13 km$^2$ (Lett and Rukhlov, 2017). Province-wide RGS compilations were released by Lett (2005), Lett (2011), and Rukhlov and Naziri (2015). Han and Rukhlov (2017) presented a data model in which all RGS data were consolidated into a unified relational database. Currently, the data set has about 65,000 samples and 5 million analyses generated by 18 analytical methods at 18 laboratories.

In the updated RGS geochemical data model (Fig. 4) we added three new entities to the skeleton data model: **data-lake** (with attributes specific to lake sediment and water samples), **data-stream** (with attributes specific to stream sediment and water samples), and **data-moss** (with attributes specific to sediment trapped by moss mats). Each of the three has a one-to-one relationship with the **data-sample** entity, meaning that theoretically, they could all be incorporated into the **data-sample** entity. However, doing so would result in many blank attributes for most of the records.

### 3.3. Data model for till geochemical data set

Till geochemical surveys typically sample subglacial tills, which are commonly considered a first derivative of bedrock (Shilts, 1993). Till has a relatively straightforward transport history that reflects the ice-flow history of an area, and the sediment geochemistry thus can be used to characterize and locate buried mineralization (Levson, 2001). Historically, till geochemical surveys conducted in British Columbia also collected other sediment types, including other till facies (ablation till, colluviated till) and glaciogenic sediments (glaciomarine, glaciofluvial), and colluvium. These deposits commonly have a more complex transport history, and identifying the source of geochemical anomalies may be difficult. The most recent regional-scale till geochemical data release (Bustard et al., 2017) was compiled from 39 reports published between 1992 and 2017 by the BCGS, GSC, and Geoscience BC. The data set has geochemical data for 10,454 samples derived from analyzing the clay (<2 μm) and silt plus clay (<63 μm) size fractions by methods including: inductively coupled plasma mass spectrometry (ICP-MS) or inductively coupled plasma emission spectrometry (ICP-ES) after an aqua regia (or modified aqua regia) digestion; lithium metaborate fusion; and instrumental neutron activation (INAA).

The reconstructed till data model (Fig. 5) fits the till geochemical data model used by Bustard et al. (2017) without adding entities or modifying relationships between entities. Only a few attributes were added to the **data-sample** and **data-analyte** entities, including *azimuth*, *dip*, *size_frac* (for size fraction), which are specific for till geochemical data.

### 3.4. Data model for coal ash chemical data set

Coal ash is the inorganic residue remain after coal combusts. It is composed of oxides of the mineral content in the coal. Coal ash chemistry can have a significant influence on coke strength after reaction (CSR), an important measure of coking coal quality. Riddell and Han (2017) designed a data model and database and filled it with coal ash oxides and related analyses for 478 samples from the Gates, Gething, Minnes and Boulder Creek formations in the Peace River coalfields in northeastern British Columbia, and from the Mist Mountain Formation in the Elk River and Crowsnest coalfields of southeastern British Columbia.

As with the till geochemical data, we found that the skeleton data model suits the coal ash chemical data well (Fig. 6). Using the same entities and relationships, the skeleton model with addition of a few coal ash specific attributes represents and describes the coal ash data completely. The additional attributes, including *coal_deposit*, *seam*, and *basis*, were added to **data_sample** and **data_analyte** entities.

129
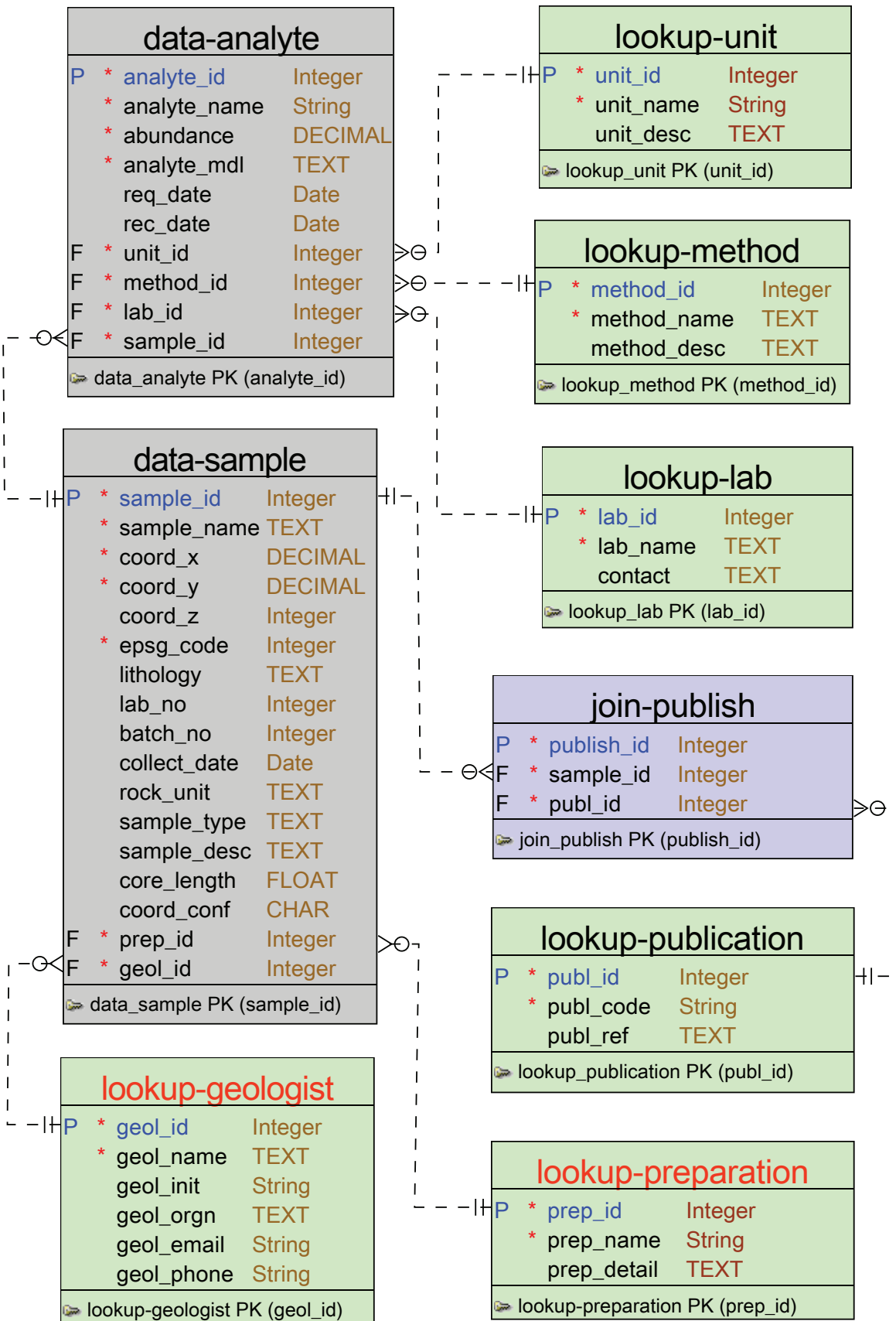
Geological Fieldwork 2018, British Columbia Ministry of Energy, Mines and Petroleum Resources, British Columbia Geological Survey Paper 2019-01

**Fig. 3.** Lithogeochemical data model reconstructed from the data model used by Han et al. (2016) by customizing the skeleton data model shown in Figure 1.
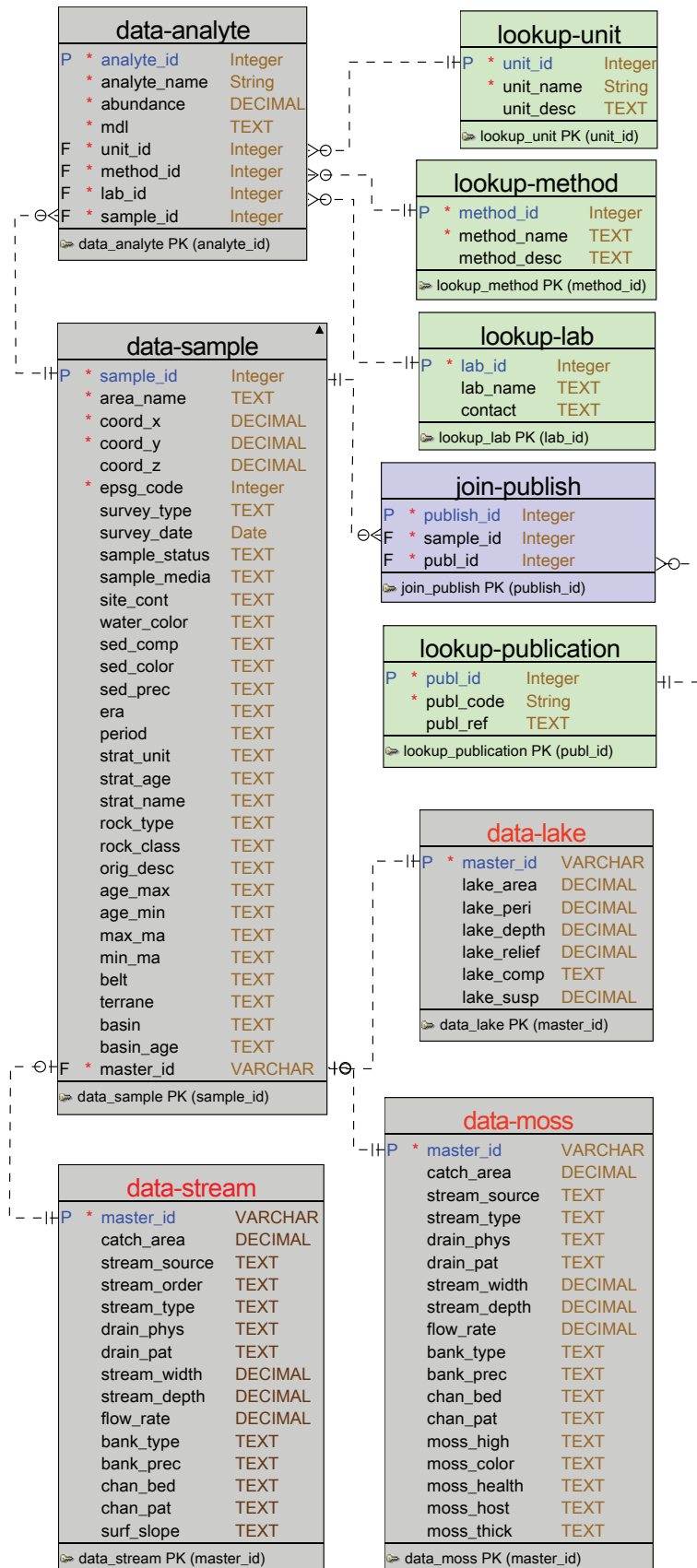
130

Geological Fieldwork 2018, British Columbia Ministry of Energy, Mines and Petroleum Resources, British Columbia Geological Survey Paper 2019-01

**Fig. 4.** RGS data model reconstructed from the data model used by Han and Rukhlov (2017) by customizing the skeleton data model shown in Figure 1.

131

Geological Fieldwork 2018, British Columbia Ministry of Energy, Mines and Petroleum Resources, British Columbia Geological Survey Paper 2019-01
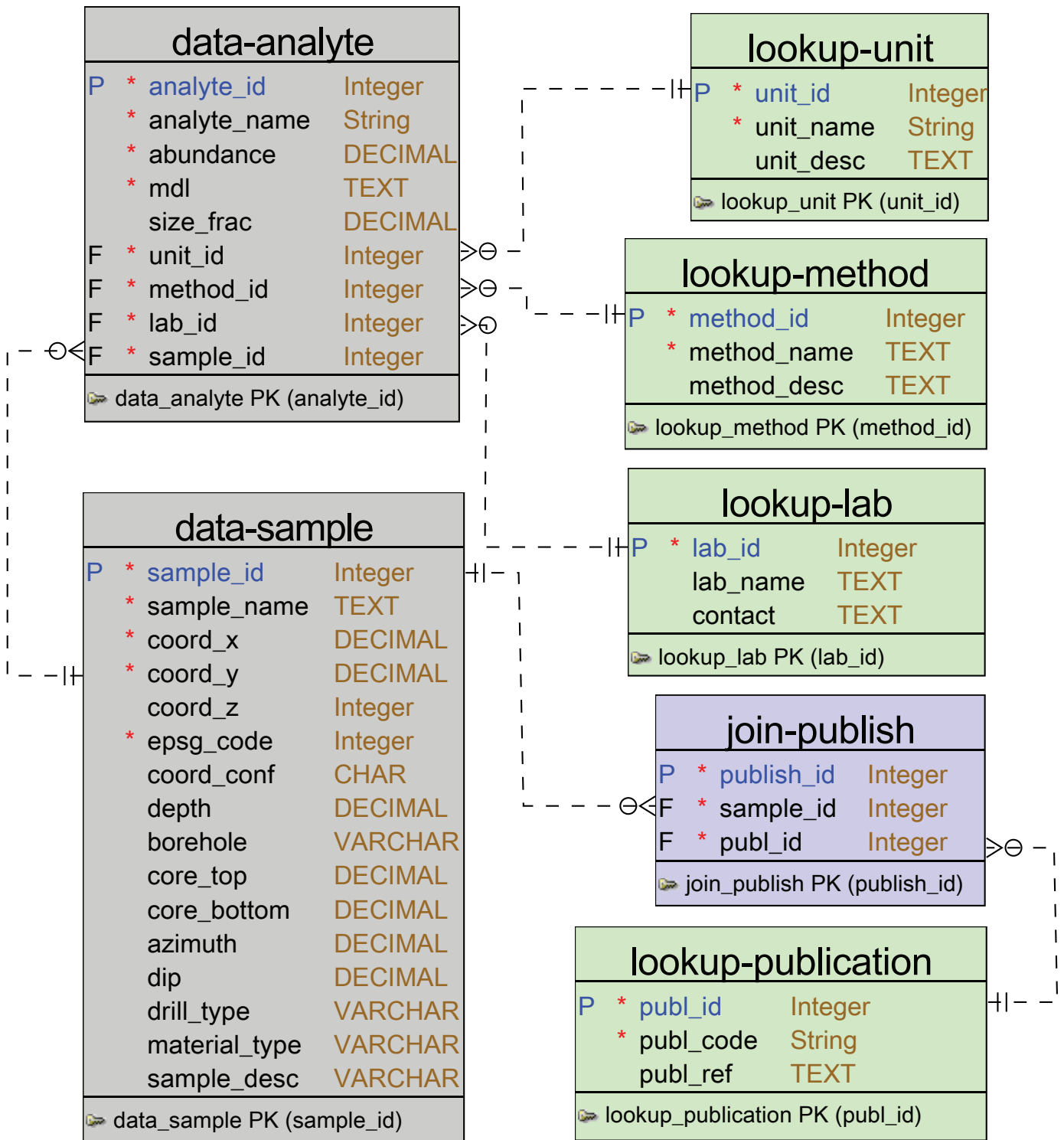
**Fig. 5.** Till geochemical data model reconstructed from the data model used by Bustard et al. (2017) by customizing the skeleton data model shown in Figure 1.
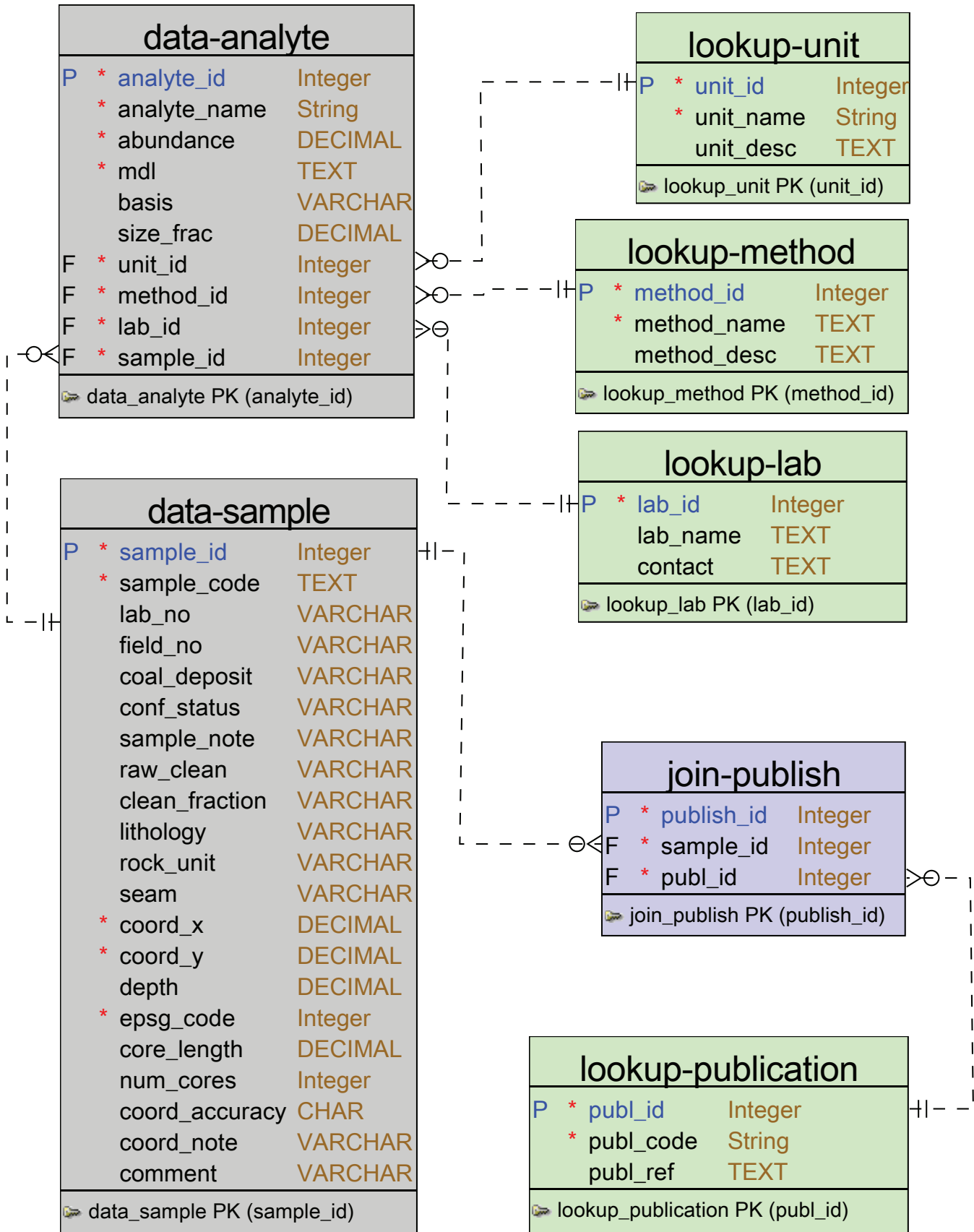
132

Geological Fieldwork 2018, British Columbia Ministry of Energy, Mines and Petroleum Resources, British Columbia Geological Survey Paper 2019-01

**Fig. 6.** Coal ash chemical data model reconstructed from the data model used by Riddell and Han (2017) by customizing the skeleton data model shown in Figure 1.

133

Geological Fieldwork 2018, British Columbia Ministry of Energy, Mines and Petroleum Resources, British Columbia Geological Survey Paper 2019-01

## 4. Development of BCGS geochemical databases

A database can be developed by implementing the corresponding data model using any database language. For the BCGS geochemical databases, we implemented the data models described above using Microsoft Access. Among the many considerations needed to develop a reliable database we paid close attention to: 1) supporting data lifecycle activities; 2) accounting for geochemical data that change with time; and 3) integrating data and metadata.

### 4.1. Supporting data lifecycle

Geoscience data, including geochemical data, live a lifecycle of six stages, from planning, acquisition, processing, analysis, storage, to publication or sharing, as defined by Faundeen et al., (2013). Except for the first two, these stages commonly occur in a database environment. This is particularly the case when geoscience data need to be long-lived, continually updated, and used by people from both within and outside an organization.

### 4.2. Accounting for geochemical data that change with time

Geochemical data sets evolve because of changes in survey techniques, improvements in data acquisition and advances in analytical instrumentation. Furthermore, previously collected and analyzed geological samples are commonly re-analyzed using different standards and techniques. For example, about 5500 of the RGS stream-sediment samples were reanalyzed using ICP-MS, with improved detection limits relative to earlier

analytical methods (Jackaman, 2017). Failing to accommodate such changes in the data model could limit the usefulness of the database.

### 4.3. Integrating data and metadata

A geochemical data set typically contains two types of data: raw data (represented by **data-sample** and **data-analyte** entities in the skeleton data model); and metadata (represented by entities with names beginning with **"lookup"** for analytical method, lab, and value unit). Results from samples with, for example, high Au contents (raw data) are meaningless if the analytical method (metadata) used to determine concentrations, are unknown. To prevent separation or loss of metadata, we unite the metadata with the corresponding raw data by enforcing the foreign key constraint between the related data and metadata entities.

## 5. Operation

To operate the four province-wide BCGS geochemical databases we built four sets of applications, using the Python scripting language, to interface with the database. These applications automate routine data management tasks, including geochemical data compilation, quality control, update, and product generation. This flow of data through these programs consists of five steps (Fig. 7): 1) data compilation; 2) initial quality control and quality assurance (QA/QC); 3) data loading; 4) product generating; and 5) product QA/QC.
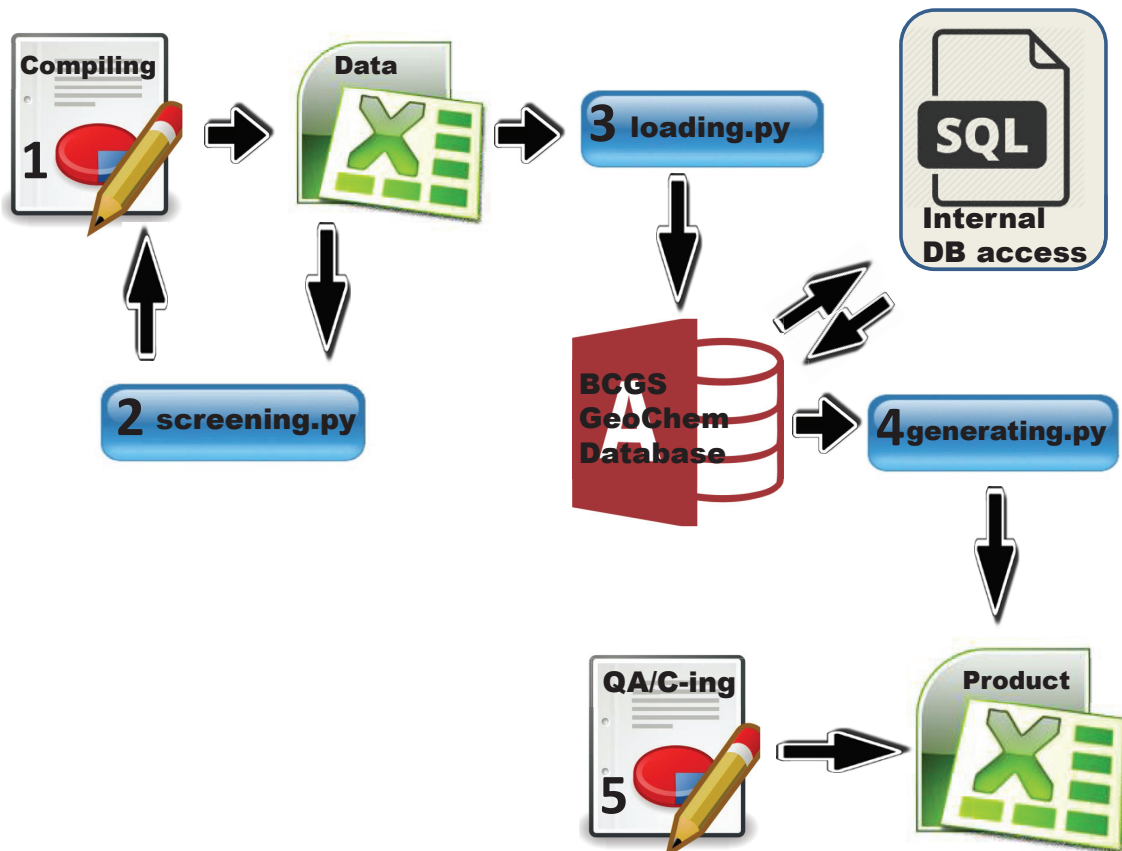


**Fig. 7.** British Columbia Geological Survey geochemical data flow.

134

Geological Fieldwork 2018, British Columbia Ministry of Energy, Mines and Petroleum Resources, British Columbia Geological Survey Paper 2019-01

During compilation (Step 1), data are retrieved from different sources and saved as Excel files in a predetermined format. Data QA/QC (Step 2) is then conducted using the corresponding Python script (screening.py), designed to flag common data errors either present in the source or introduced during data compilation. The errors flagged may include unrealistic determinations and units, improper methods, inconsistent handling of censored data (e.g., values below the detection limit), wrong sample locations, and redundant samples. The flagged errors are then manually examined and corrected. After this step, data are loaded into the database (Step 3). This is done automatically by executing the Python script (loading.py). Generating derived data products (Step 4) is also done using a Python script (generating.py), which retrieves and outputs data in simple formats, such as Comma Separated-Value (CSV), ESRI shapefiles, or MS Excel files. If errors are found in the generated data products, Steps 1 to 4 are repeated.

The geochemical databases discussed in this paper are not designed for direct access by data users but for data management personnel who are responsible for operating, maintaining, and updating these databases. As indicated in Figure 7, data management personnel prepare and release data products derived from these databases to users in simple tabular formats.

## 6. Conclusion

In this paper we present a simplified geochemical data model that is capable of capturing and representing generic entities, common attributes, and intrinsic relationships existing across different geochemical data sets. This skeleton data model has the potential for us to consolidate the four province-wide geochemical databases that currently operate independently into a unified one, improving the efficiency and standardization of our geochemical data management.

Data modeling is typically an incremental process. It is common to start with a simple data model that satisfies immediate needs and to later add in complexities to meet requirements that were unforeseen in the initial analysis of database requirements. The skeleton data model is not meant to be one that is all-inclusive, all-purpose. For example, because the geochemical data currently stored in the BCGS geochemical databases are only from field samples, we kept the data model simple and excluded analytical duplicates, blanks, and reference materials. But the skeleton data model is capable of expanding to include such samples. It currently includes all the basic elements found across geochemical data sets at the BCGS, and can be built on.

## Acknowledgments

## References cited

Bustard, A.L., Han, T., and Ferbey, T., 2017. Compiled till geochemical data for British Columbia. British Columbia Ministry of Energy and Mines, British Columbia Geological Survey GeoFile 2017-9, 7 p.

Connolly, T., and Begg, C., 1999. Database Sytems – A Practical Approach to Design, Implementation, and Management, 2nd Edition, Addison-Wesley, 113 p.

International Association of Oil & Gas Producers (IOGP), EPSG Geodetic Parameter Dataset, http://www.epsg-registry.org/, 2018.

Faundeen, J.L., Burley, T.E., Carlino, J.A., Govoni, D.L., Henkel, H.S., Holl, S.L., Hutchison, V.B., Martín, E., Montgomery, E.T., Ladino, C.C., Tessler, S., and Zolly, L.S., 2013. The United States Geological Survey Science Data Lifecycle Model: U.S. Geological Survey Open-File Report 2013-1265, 4 p., http://dx.doi.org/10.3133/ofr20131265.

Granitto, M., Schmidt, J.M., Labay, K.A., Shew, N.B., and Gamble, B.M., 2012. Alaska geochemical database-Mineral exploration tool for the 21st century. U.S. Geological Survey Open-File Report 2012-1060, 33 p.

Han, T., Rukhlov, A.S., Naziri, M., and Moy, A., 2016. New British Columbia lithogeochemical database: Development and preliminary data release. British Columbia Ministry of Energy and Mines, British Columbia Geological Survey GeoFile 2016-4, 6 p.

Han, T., and Rukhlov, A.S., 2017. Regional Geochemical Survey (RGS) data update and release using the newly developed RGS database. British Columbia Ministry of Energy and Mines, British Columbia Geological Survey GeoFile 2017-11, 7 p.

Jackaman, W., 2017. RGS Sample reanalysis project (parts of NTS 082G, 082J, 092N, 093E, 093H, 103O, 103P and 104N), GeoScience BC Report 2017-04.

Lett, R.E., 2005. Regional Geochemical Survey Database on CD. British Columbia Ministry of Energy and Mines, British Columbia Geological Survey GeoFile 2005-17.

Lett, R.E., 2011. Regional Geochemical Survey Database 2011. British Columbia Ministry of Energy and Mines, British Columbia Geological Survey GeoFile 2011-7.

Lett, R., and Rukhlov, A.S., 2017. A review of analytical methods for regional geochemical survey (RGS) programs in the Canadian Cordillera. In: Ferbey, T., Plouffe, A., and Hickin, A.S., (Eds.), Indicator Minerals in Till and Stream Sediments of the Canadian Cordillera. Geological Association of Canada Special Paper Volume 50, and Mineralogical Association of Canada Topics in Mineral Sciences Volume 47, pp. 53-108.

Levson, V.M., 2001. Regional till geochemical surveys in the Canadian Cordillera: sample media, methods, and anomaly evalutation. In: McClenaghan, M.B., Bobrowsky, P.T., Hall, G.E.M., and Cook, S.J. (Eds.), Drift Exploration in Glaciated Terrain, Geological Society of London, Special Publication 185, pp. 45-68.

Riddell, J., and Han, T., 2017. Ash chemistry database for British Columbia Rocky Mountain bituminous coals. Ministry of Energy and Mines, British Columbia Geological Survey GeoFile 2017-10, 15 p.

Rukhlov, A., and Naziri, M., 2015. Regional Geochemical Survey database. British Columbia Ministry of Energy and Mines, British Columbia Geological Survey GeoFile 2015-3.

Shilts, W.W., 1993. Geological Survey of Canada's contributions to understanding the composition of glacial sediments. Canadian Journal of Earth Sciences, 30, 333-353.

Watson, C., and Evans, I., 2012. Geochemistry data model summary, British Geological Survey, http://www.bgs.ac.uk/services/dataModels/geochemistry.html. Last accessed November 2016.