# Knowledge graphs for seismic data and metadata

William Davis [a,b,*], Cassandra R. Hunt [b]

[a] Cecil H. and Ida M. Green Institute of Geophysics and Planetary Physics, Scripps Institution of Oceanography, University of California, San Diego, La Jolla, CA 92037, USA
[b] RelationalAI, Berkeley, CA, USA

## ARTICLE INFO

## ABSTRACT

The increasing scale and diversity of seismic data, and the growing role of big data in seismology, has raised interest in methods to make data exploration more accessible. This paper presents the use of knowledge graphs (KGs) for representing seismic data and metadata to improve data exploration and analysis, focusing on usability, flexibility, and extensibility. Using constraints derived from domain knowledge in seismology, we define a semantic model of seismic station and event information used to construct the KGs. Our approach utilizes the capability of KGs to integrate data across many sources and diverse schema formats. We use schema-diverse, real-world seismic data to construct KGs with millions of nodes, and illustrate potential applications with three big-data examples. Our findings demonstrate the potential of KGs to enhance the efficiency and efficacy of seismological workflows in research and beyond, indicating a promising interdisciplinary future for this technology.

## 1. Introduction

Navigating big data is becoming increasingly crucial for seismic studies of the Earth's structure, tectonic processes, and related geohazards (Arrowsmith et al., 2022). Collectively the field of seismology generates vast amounts of diverse data in many formats, including time-series waveforms, metadata pertinent to the instruments and stations which record them, and catalogues of estimated event source parameters. For instance, the Incorporated Research Institutions for Seismology (IRIS) Data Management Center (DMC) provides access to over 850 TB of archive data, including waveform, station, and event metadata across more than 27 data formats, as well as other higher-level data products (Trabant et al., 2012; Hutko et al., 2017). The scale and diversity of data sources and schema complicate data exploration, especially where sifting through large volumes or joining across sources is required (Dost et al., 2009; Krischer et al., 2016; Ringler et al., 2022; Arrais et al., 2022). Nuanced data requirements result in bottlenecks where researchers first bulk download records and then refine using hand-crafted data transformation and analysis workflows. Effective data utilization is further challenged by the rapid acceleration of data generation, primarily driven by the development of new data-dense, distributed sensor systems (Zhan, 2020; Trugman et al., 2022; Spica et al., 2023). Traditional methods of utilizing these data rely on specialized software tools, dataframe analysis libraries (e.g., Pandas), and database systems, requiring researchers

to navigate complicated schema outlines or data format specifications. There is increasing recognition that seismic data must be made more accessible, both to improve the research pipelines of the research seismological community (Gil et al., 2018; Arrowsmith et al., 2022), but also to facilitate broader applications to geohazard assessment, oil and gas exploration, data science, and machine learning domains (Mohammadpoor and Torabi, 2020; United States Geological Survey (USGS) (2021); Ringler et al., 2022). To serve diverse end goals, seismic data exploration must be flexible and accessible. As new data sources become available, exploration methods must be extensible to accommodate them.

One route to improve data accessibility utilizes graphical user interface-based web services (e.g., Weertman, 2010; Newman et al., 2013; Falco et al., 2017). These tools enable access to homogeneous data through a single interface, allowing users to query seismic data, for example, based on event parameters—such as location, time, and magnitude. However, these tools are in practice restricted to specific data sources and data search is simplified in a way that restricts query complexity. Recently, Yu et al. (2021) used cloud-based services to offer a route to scalable storage and computation for seismic data access and analysis. The catalog, hosted by the Amazon Web Services (AWS) Open Dataset Program initiative, brings multiple data sources from the Southern California Seismic Network (SCSN) together in a single "data lake". The records, stored on an AWS bucket, are searchable via metadata in the names of files or filtering on certain data values recorded in
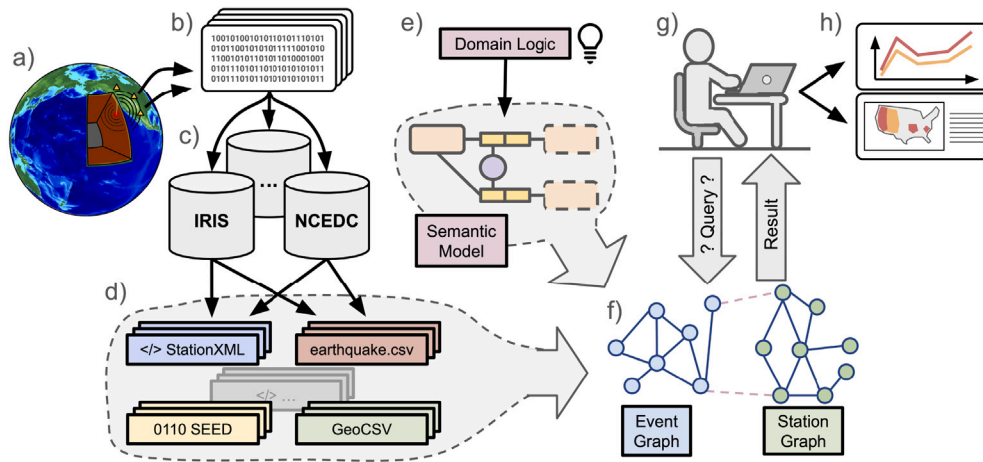
**Fig. 1.** A visual schematic of our approach to knowledge graphs (KGs) in a seismic data workflow. (a) Ground motion from earthquakes or other sources is recorded by seismometers or other instruments. (b) Raw instrument data is collected, stored, transformed, and managed by (c) various seismic data centers and facilities. (d) Higher-level data files and data products, such as earthquake catalogs, are produced and made available by the data management facilities. (e) Domain knowledge is used to create a semantic model for the KGs. (f) A KG database is populated from source data using logic derived from the semantic model. (g) The user queries the KGs. (h) The queried data is use for science goals .

index files using the AWS Command Line Interface (Southern California Earthquake Data Center (SCEDC), 2021).

An alternate and potentially complementary approach is to map heterogeneous data schemas to a common, extensible, and queryable semantic model. Data integration using a common ontology may be realized virtually, with mediated approaches (Halevy et al., 2006; Xiao et al., 2019), or physically in a single database. Recently, knowledge graphs (KGs) have emerged as a promising approach to organize complex and interconnected data in ontologies (Hogan et al., 2021; Gutiérrez and Sequeda, 2021), which can be tailored to meet specific requirements and domains (Abu-Salih, 2021). KGs are being increasingly utilized in geosciences (see Ma, 2022, for a comprehensive review). The use of KGs offer a versatile and extensible solution for many aspects of the data life-cycle, from data representation and curation, integration, and data analysis and result communication (Ma et al., 2014; Wing, 2019).

This paper introduces the idea of using relational KGs for seismic data, delivering a queryable semantic model and addressing the challenges in data exploration with large and schema-diverse seismic data. In this way, KGs complement web service and data lake offerings. We emphasize two key benefits of representing geoscience data with KGs: (1) scalability and performance competitive with modern SQL databases (Timón-Reina et al., 2021; Monteiro et al., 2023; Hölsch et al., 2017), and (2) ability to combine structured and semistructured source data in a common representation, extensible to new data attributes and sources. These advantages render KGs uniquely amenable to the evolving data landscape of seismology. We first outline a semantic model consisting of two KG ontologies, one for seismic station metadata and one for earthquake event data. We then present an implementation of these KGs demonstrating the integration of 4 data sources into a common, searchable graph structure, and provide three example applications. Our approach is diagrammed in Fig. 1. Our KGs are constructed from declarative definitions, enabling the abstraction of implementation details and a focus on knowledge modeling (Humphries, 2021). The KG definitions utilize a physical data integration approach, with definitions materialized on-demand, taking advantage of a recently developed scalable, cloud native relational KG management system (RKGS). We emphasize that we are not introducing a new data format; we are introducing KGs as a "semantic layer" for seismic knowledge (Stirewalt and Búr, 2023), to augment and connect heterogeneous data from existing sources.

## 2. KGs for Seismic knowledge

In this study, we model two types of seismic knowledge: station metadata and seismic event data. In seismology, station metadata denotes known information about seismic stations and seismometers, such as geographic location, orientation, local site effects, and instrument type. Conversely, event data, gathered in earthquake catalogs such as the Global Centroid-Moment-Tensor (GCMT) project (Dziewonski et al., 1981; Ekström et al., 2012), describes earthquakes and other anthropogenic activities by their estimated properties, such as location, moment magnitude, and depth. This data differs from station metadata as it is based on inferences of natural events, involving uncertain, idealized representations of physical phenomena. Another type of seismic data is waveform data generated by seismic instruments, however, for simplicity, we choose not to include this in our current study.

We represent seismic knowledge in a graph structure. Nodes represent abstract objects (e.g., *the Berkeley Digital Seismic Network* or *the Columbia College Station*). Nodes can also represent atomic property values, like a specific latitude (e.g., *37.9°*). Edges describe relations between objects (e.g., the Berkeley Digital Seismic Network *manages* the Columbia College Station). An example KG is shown in Fig. 2.(a).

The nodes and edges in a KG organize data according to an ontology: a formal description of the concepts and relationships within a domain. We diagram the ontologies of seismic knowledge with Object-Role Modeling (ORM) (Halpin, 2015). We choose ORM to represent each ontology as it captures the relationship between nodes and edges as well as data constraints important to populating the KG, as we will show later. Importantly, ORM is attribute-free, modeling all relationships as explicit facts and disentangling ontology semantics from a specific KG implementation. The ontology is applied here to build a relational KG, but may be equally applied to a labeled property graph (LPG) or Resource Description Framework (RDF) graph, for example. An example ORM diagram, without data constraints, is shown in Fig. 2.(b).

To model seismic knowledge as relational KGs, we define KG ontologies through the recognition of data integrity constraints, declared in natural language. These directly correspond to fact types in the ontology diagram and determine the relevant entity, value, and edge relations, including the criteria for uniquely identifying each entity. In the following section, we propose ontologies for station metadata and seismic event data.
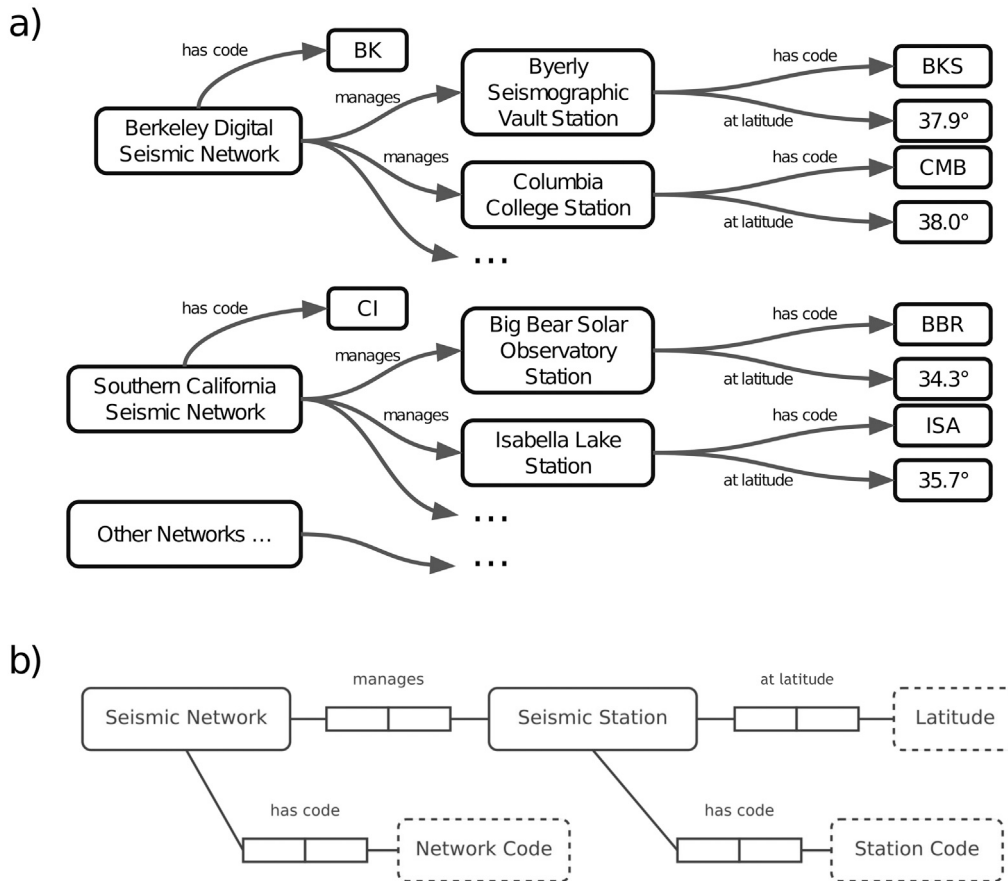
**Fig. 2.** Example KG and Object-Role Modeling (ORM) diagrams for a model of station metadata. For this illustrative example, details have been substantially simplified. Subplot (a): An example KG for a subset of station metadata. Nodes represent abstract objects (e.g., *the Berkeley Digital Seismic Network* or *the Columbia College Station*) and also atomic values with their data type (e.g., the latitude *37.9°*) and are diagrammed here using rounded boxes. In this example, the latter node type captures attribute information in the source data and is equivalent to a node property in a property graph representation, but we need not make that distinction in a relational KG. Edges describe relations between nodes (e.g., the Berkeley Digital Seismic Network *manages* the Columbia College Station) and are diagrammed using arrows. Subplot (b): An ORM diagram for the above KG for station knowledge. Nodes that represent abstract objects are labeled as an "entity type" (e.g., the Columbia College Station is a *Seismic Station*) and are diagrammed using solid-edged rounded boxes. Nodes that are self-identified by their atomic data value are labeled as a "value type" (e.g. 37.9° is a *Latitude*) and are diagrammed using dashed-edged rounded boxes. Edge labels are represented with binary "roleboxes", one connected with a line to each entity type, or to an entity type and value type. A set of roleboxes, and the entity types and value types connected to them, are referred to as a "fact type" in the ontology. For more details on ORM diagrams, see Halpin (2015).

## 2.1. Modeling station knowledge

The first type of knowledge we consider describes seismic instruments, and their hierarchical groupings and associations. We begin by identifying and verbalizing facts and constraints (S1–19) in the ontology, diagrammed in Fig. 3.

First, we identify four entities:

- `Channel`: An individual seismic instrument or sensor.
- `Channel Group`: A group of multiple channels. For practical purposes, channels are often grouped together into orthogonal triples.
- `Station`: A location—for example, a building—housing seismic instrument(s).
- `Network`: A collection of seismic stations, which are either managed and maintained by a specific agency or are linked to a specific scientific campaign.

Often, "(seismic) station" is used as a signifier for this entire hierarchy. The semantic model draws inspiration from the International Federation of Digital Seismograph Networks (FDSN) Source Identifiers specification (Trabant et al., 2019; Benson et al., 2019), and the FDSN Station Extended Markup Language (StationXML) format (see Data and Code Availability). However, we introduce augmentations to give

added utility to the model. In particular, `Channel Group` is not represented as an element in the StationXML format. We emphasize that the semantic concepts here are general, and may be mapped to station metadata represented with other schemas (e.g., Ahern et al., 2009; Schorlemmer et al., 2011).

Identifiers for each entity type node must be graph-unique. This requirement distinguishes a relational KG representation of edges and nodes from 6th normal form (6NF) (Date, 2006): each node and edge relation (or table) cannot be normalized further, as required by 6NF, and additionally the primary and foreign keys (node identifiers) must uniquely represent the same nodes across the entire set of relations in the graph. The combined requirement of 6NF representation and graph-unique node identifiers is known as "Graph Normal Form" (Stirewalt and Búr, 2023). To define the combination of data that constitutes a graph-unique identifier for each entity type, we recognize certain edge relations and integrity constrains on those relations:

- S1. **Each** `Station` is managed by **exactly one** `Network`,
- S2. **Each** `Channel Group` is in **exactly one** `Station`, and
- S3. **Each** `Channel` is in **exactly one** `Channel Group`.

Organizational bodies regularly define identification codes for networks, stations, channel groups, and channels (e.g., Buland, 2012; International Seismological Centre (ISC) (2020)). Expressed as a modeling decision, this corresponds to each entity having exactly one identification code as part of its reference scheme. We recognize that:

S4. **Each** `Network` has a code of **exactly one** network code,
S5. **Each** `Station` has a code of **exactly one** station code,
S6. **Each** `Channel Group` has a code of **exactly one** location code, and
S7. **Each** `Channel` has a code of **exactly one** channel code.

In addition to identification codes, the entities have other associated properties. Some of these properties are explicitly represented in the FDSN StationXML format. For example, we incorporate information on geographic location in our ontology, which are modeled as mandatory and single-value relations:

S8. **Each** `Station` is at **exactly one** latitude,
S9. **Each** `Station` is at **exactly one** longitude, and
S10. **Each** `Station` is at **exactly one** elevation.

Other properties define aspects of channel instrumentation and digitization. The "band type" defines the general sampling rate and response band of the data source. The "instrument type" (or "source") defines the type of sensor or data source (e.g., seismometer, accelerometer, geophone). The "orientation" (or "subsource") indicates the orientation of the measurement. The traditionally used orientations are North (N), East (E), and Up (Z). These properties are modeled as mandatory and single-valued:

S11. **Each** `Channel` has **exactly one** band type,
S12. **Each** `Channel` has **exactly one** instrument type, and
S13. **Each** `Channel` has **exactly one** orientation.

Additional properties define depth and operational extent:

S14. **Each** `Channel` is at a depth of **exactly one** depth,
S15. **Each** `Channel` was operational from **exactly one** date-time, and
S16. **Each** `Channel` is operational until **exactly one** date-time.

Finally, we define the minimum combination of data that constitutes a graph-unique preferred identifier for each entity type. We choose to encode the rules of the FDSN Source Identifiers as uniqueness constraints. `Networks` are uniquely defined by their `Network` codes (S4) (Buland, 2012; International Seismological Centre (ISC) (2020)). For the remaining entity types, uniqueness is defined by the hierarchical constraints S1–3 combined with the entity's own identification code (S5–7):

S17. **For each** `Network` **and** station code,

  • **at most one** `Station` is managed by **that** `Network` **and** has that station code.

S18. **For each** `Station` **and** location code,

  • **at most one** `Channel Group` is in **that** `Station` **and** has that location code.

As the FDSN Source Identifier specifications do not prescribe uniqueness conditions for channels—codes instead indicate instrumentation details—we choose to define the following criterion:

S19. **For each** `Channel Group` **and** channel code operational from **that** date-time,

  • **at most one** `Channel` is in **that** `Channel Group` **and** has that channel code **and** was operational from that date-time.

The start date-time requirement naturally allows enforcement of constraints S15 and S16: a `Channel` that has multiple operational periods will be represented by multiple `Channel` nodes, one for each period.

*2.2. Modeling seismic event knowledge*

We now model knowledge associated with records of seismic events in catalogs. As this knowledge reflects idealizations of natural events, records of the same natural event may vary in both schema and data, which the structure of an ontological model should handle. We identify facts and constraints (E1–10) that promote an event knowledge model flexible enough to encompass data from many sources, diagrammed in Fig. 4.

We define two entities associated with event knowledge:

• `Contributor`: An agency or group that manages, maintains, and contributes data to a seismic event catalog.
• `Event Record`: A record or entry of a seismic event.

We use the term "(seismic) event" as a signifier of this ontology. We define a mandatory and single-valued binary relation between these entities:

E1. **Each** `Event Record` is contributed by **exactly one** `Contributor`.

Note that we model the concept of an `Event Record` in a catalog rather than attempting to model the physical event itself. If one earthquake appears in two catalogs, our model will regard them as two independent event records. Subsequent entity resolution—or deduplication—may be used to associate event records with a unique seismic event (Sun et al., 2020; Obraczka et al., 2021). Each entity has a mandatory and single-valued reference scheme, which we verbalize as:

E2. **Each** `Contributor` has a code of **exactly one** contributor code, and
E3. **Each** `Event Record` has **exactly one** event ID.

We also model properties of the `Contributor` and `Event Record` entities. For each `Contributor`, we include a mandatory (but not necessarily single-valued) catalog code:

E4. **Each** `Contributor` has a catalog code that is **some** catalog code.

In the GCMT catalog, for example, this refers to the "hypocenter reference catalog" code. Each `Event Record` has associated property values corresponding to estimated physical parameters of the event. We choose to incorporate a small but fundamental set of (mandatory and single-valued) properties in our ontology:

E5. **Each** `Event Record` has a magnitude of **exactly one** magnitude,
E6. **Each** `Event Record` occurred at **exactly one** origin date-time,
E7. **Each** `Event Record` was at a latitude of **exactly one** latitude,
E8. **Each** `Event Record` was at a longitude of **exactly one** longitude, and
E9. **Each** `Event Record` was at a depth of **exactly one** depth.

Finally, we define the graph-unique preferred identifiers for each entity type. By constraint E2, a `Contributor` is uniquely defined by their contributor name. For an `Event Record`, we require that event IDs are unique within catalogs. This is modeled as an external uniqueness constraint over relations E1 and E3:

E10. **For each** `Contributor` **and** event ID,

  • **At most one** `Event Record` was contributed by **that** `Contributor` **and** has that event ID.
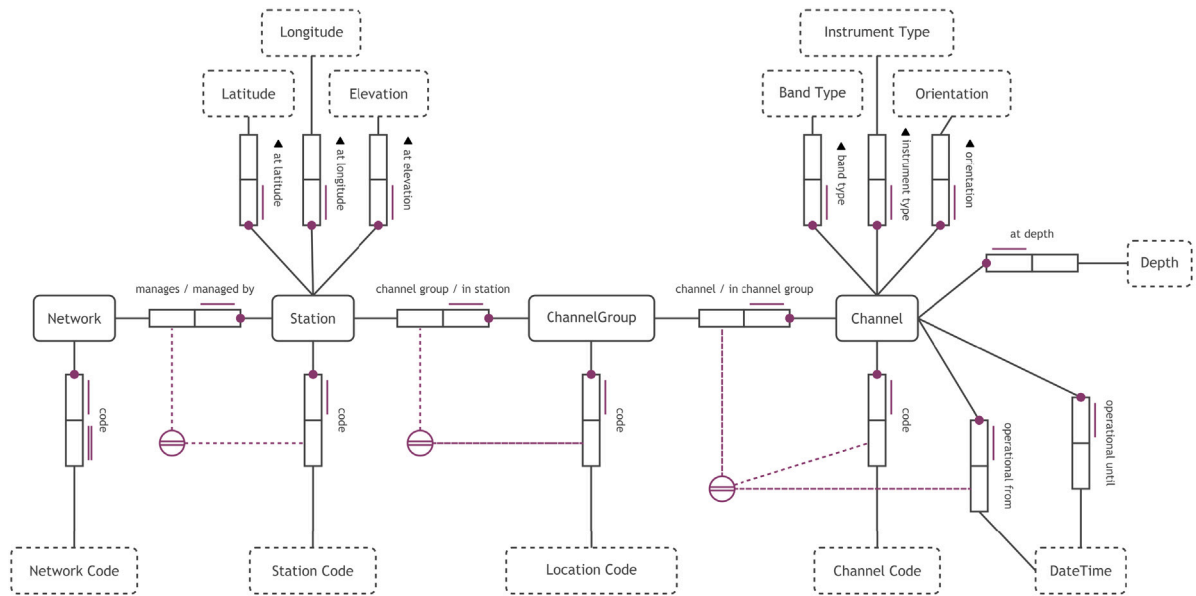
**Fig. 3.** Our ORM diagram for the station knowledge ontology. Entity types (e.g., Station) are represented by solid-edged rectangles. Value types (e.g., Elevation) are represented by dashed-edged rectangles. Binary fact types—for example, S1: "**Each** Station is managed by **exactly one** Network"—are represented by entity and value types connected to a pair of roleboxes, along with a set of constraints (in violet). Edge names are indicated with text next to the roleboxes. Violet lines next to the roleboxes indicate uniqueness constraints, whereas violet dots indicate mandatory roles. A double violet line indicates a preferred identification scheme, for example, constraint S4. Violet dashed lines leading to violet ⊜ symbols correspond to an external preferred identification scheme, for example, constraints S17–19.
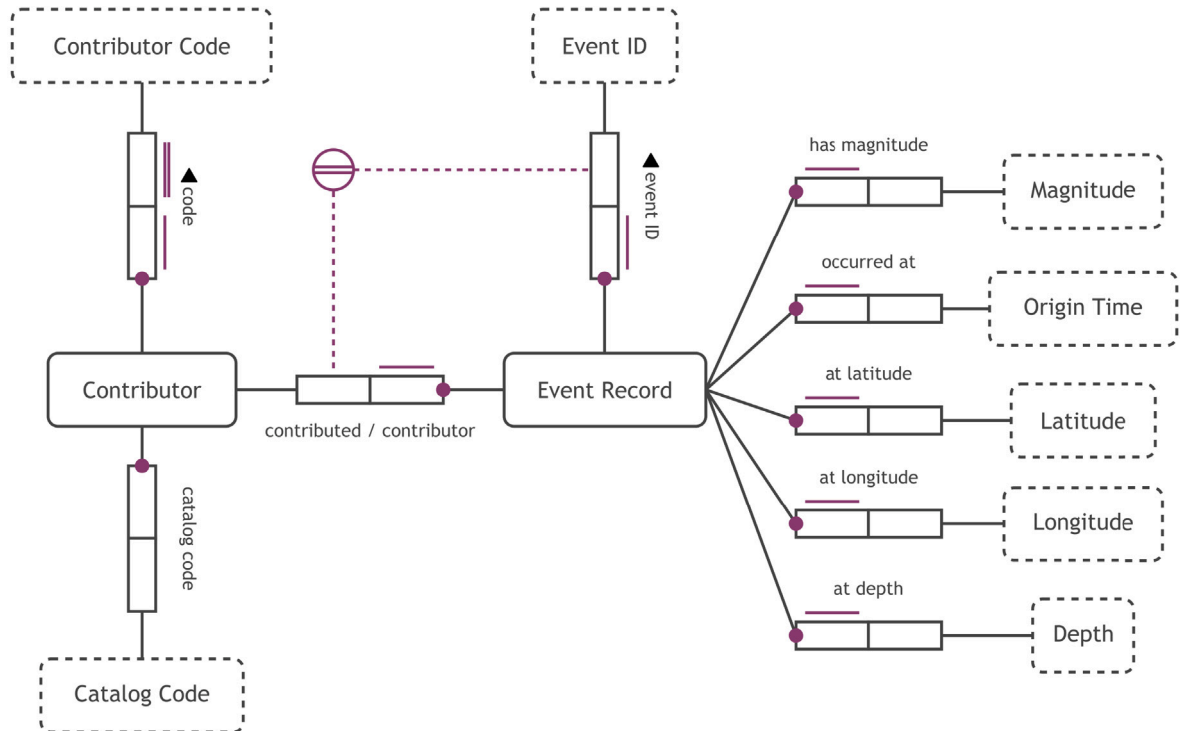


**Fig. 4.** Our ORM diagram for the event knowledge ontology. Contributor and Event Record entity types are represented by solid-edged rectangles. Each type of property value (e.g., Depth) is represented by a dashed-edged rectangle. Binary fact types—for example, E9: "**Each** Event Record has a depth of **exactly one** depth".—are represented by entity and value types connected to a pair of roleboxes, along with a set of constraints (in violet). Edge names are indicated with text next to the roleboxes. Violet lines next to the role boxes indicate uniqueness constraints, whereas violet dots indicate mandatory roles. The double violet lines for contributor name signifies it as a preferred identification scheme. The violet dashed lines leading to the violet ⊜ symbol corresponds to the external preferred identification scheme in constraint E10.

## 3. Implementation of station and event KGs

To study the functionality of the two proposed KGs, we develop an implementation of the station and event ontologies in a database. This is accomplished using the RelationalAI RKGS (RelationalAI (RAI),

2021b), and modeled using the declarative, relational language Rel (RelationalAI (RAI) (2021a); Stirewalt, 2022). We de-emphasize language-specific details in favor of providing an outline of the process of mapping seismic data into KGs (all code is available in the supplementary material). With the ontology of our two KGs outlined in

the previous section, we now focus on populating the graphs with real-world seismic data (Hofer et al., 2023).

### 3.1. Data selection and extraction

We identify a range of relevant sources of seismic data to integrate into our KGs. These sources highlight the data-schema diversity present in file formats commonly used by seismologists. We consider:

- Station metadata, in StationXML format, acquired from IRIS DMC using the `fdsnws-station` webservice (see Data and Code Availability),
- Earthquake event data, in NDK format, acquired from the Global Centroid-Moment Tensor (GCMT) catalog webservice (Dziewonski et al., 1981; Ekström et al., 2012),
- Earthquake event data, in CSV format, acquired from the Northern California Seismic Network (NCSN) catalog using the NCEDC's Northern California Earthquake Catalog Search webservice (NCEDC, 2014), and
- Earthquake event data, in CSV format, acquired from the United States Geological Survey (USGS) earthquake catalog webservice (United States Geological Survey (USGS) and Earthquake Hazards Program, 2017).

Where multiple event data are available, we use the most recent, preferred solution. The precise search parameters for extracting data from these sources vary depending on the intended application of the KG and are specified in Section 4.

### 3.2. Data loading and transformation

We employ rule-based, declarative relation definitions, written in the Rel language, to transform both structured and semi-structured data sources to a relational KG. With this approach, the transformation logic, source data, and KG may coexist in the same database, preserving data provenance and allowing queries across graph and source data. The transformation logic takes advantage of Rel's support for entity generation, querying over schema, higher order logic, and data integrity constraint declarations. However we note that the extract-load-transform process need not be constrained to one approach for all data sources. For example, data transformation between structured formats using domain specific languages has been widely studied (García-González et al., 2020; Hofer et al., 2023).

Mapping input data to KG values requires knowledge of the schema for each data format. For example, in NDK format a magnitude estimate is located in the character range 49–55, whereas NCEDC CSV data stores equivalent information in the "Magnitude" column. In another example, the band type, instrument type, and orientation of a `Channel` can be inferred from the channel code using the FDSN Source Identifiers. Our implementation populates the station KG with `Network`, `Station`, and `Channel` property labels and values from StationXML data by querying over the source data schema.

Population of edge relations between entity types also differs for each data format. The hierarchical structure of StationXML enables the relations S1–3 to be inferred directly from attributes and sub-elements outlined in the StationXML specification. For the event data, the tabular structure of the source data allows relation E2 to be realized by identifying data appearing in a common row, (or, for NDK files, sets of rows).

### 3.3. Entity creation

Entity identifiers are represented as hashes of their node label plus the preferred identification data which uniquely identify each node, as declared in Section 2. In the station KG, uniqueness is identified for a `Network` from the extracted network code (S4). For the remaining station graph entities, we invoke the external uniqueness constraints S17–19, defining:

- Each `Station` node identifier as a hash of `Network` node identifier and the station code,
- Each `Channel Group` node identifier as a hash of the `Station` node identifier and the location code, and
- Each `Channel` node identifier as a hash of the `Channel Group` node identifier, channel code, and start date-time.

For the event KG, uniqueness for `Contributor` entities is identified through the extracted name (E2). With `Contributor` entities resolved, external uniqueness constraint E10 is invoked, such that:

- Each `Event Record` node identifier is a hash of its Event ID and `Contributor` node identifier.

### 3.4. Quality assurance

For the KGs to be trustworthy and useful for analysis applications, the correctness of the mapped knowledge must be verified (Wang and Strong, 1996). We codify data constraints S1–19 and E1–10 into logical rules—or programmatic integrity constraints—to identify aberrations or logical errors in the KG that may have arisen during construction. If any integrity constraint is violated, the construction of the KG will halt, and the data can be interrogated for aberrations. We note that the data collected across all sources in this study were typically of high quality. Only one case required data cleansing: a StationXML file contained a duplicate channel which was manually removed.

## 4. Querying examples

To empirically demonstrate the effectiveness of seismic KGs, we show three examples using data described in Section 3.1. These examples range from very simple queries that are readily achievable using existing tools (e.g., Weertman, 2010; Beyreuther et al., 2010; Newman et al., 2013), to more complex queries that take advantage of the KG structure and logical rules.

### 4.1. Example 1: Filtering events by location

We first investigate executing simple queries and determining aggregate measures using a single KG. The objective is to filter a large set of event data—from multiple sources—by geographic position. This example demonstrates data integration of a mix of semi-structured and tabular source data.

We collect all available event data for the year 2020 from the three event sources listed in 3.1 and construct an event KG, `event_kg`. The resulting KG comprises over 1.6 million nodes and 2.3 million edges in the ontology, including 18 `Contributor` entities and approximately 230,000 `Event Record` entities. We filter the events by geographic position around California and calculate the number of contributions from each contributor. Code for this query is given in Listing 1, and the resulting aggregated counts are shown in the legend of Fig. 5. We take advantage of Rel's support for ungrounded relations to define reusable logic `filter_latitude` and `filter_longitude`, which are evaluated on demand in the query.

To further interrogate the event data, we query the KG to determine the spatial distribution of event epicenters for each contributor, and output the table `list_positions`. Code for this statement is given in Listing 2, and the output table is used to generate Fig. 5.

The overlapping points in Fig. 5 indicate that many `Event Record` entities from different contributors may correspond to the same physical events. This could arise from catalogs sharing a common data origin, as exemplified between contributors NCSN (acquired from the NCEDC) and NC (acquired from the USGS). Identifying these matches would require entity resolution (Sun et al., 2020; Obraczka et al., 2021).

Listing 1: Querying the event KG (`event_kg`) through the constraints outlined in Example 1. Lines 1–3 filter the event KG by latitude and longitude (definitions of `filter_latitude` and `filter_longitude` are given in the supplemental material). Lines 5–9 define a binary relation of contributor names and filtered event nodes by performing an inner join over `contributed` and `name`. Lines 11-12 define a binary relation of contributor names and the total number of events contributed.

```
1   def filter_event_CA(event) =
2       filter_latitude(event_kg,32.6,42.6,event) and
3       filter_longitude(event_kg,-126.2,-113.7,event)
4
5   def filtered_contributor_event(contributor_name,event) =
6       filter_event_CA(event) and
7       event_kg:contributed(contributor,event) and
8       event_kg:name(contributor,contributor_name)
9       from contributor in event_kg:Contributor
10
11  def count_events(contributor_name,event_total) =
12      count(filtered_contributor_event[contributor_name],event_total)
```

Listing 2: Querying the event KG (`event_kg`) through the constraints outlined in Example 1. The code defines a tabular view of contributor names and filtered latitudes and longitudes by performing an inner join over `filtered_contributor_event`, `at_latitude`, and `at_longitude`. The format is column, row, value, with the `event` node serving as a unique row identifier.

```
1   def list_positions =
2       :Name, event, contributor_name;
3       :Latitude, event, latitude;
4       :Longitude, event, longitude
5       from event in event_kg:EventRecord, contributor_name, latitude, longitude
6       where
7           filtered_contributor_event(contributor_name,event) and
8           ^Latitude(latitude, event_kg:at_latitude[event]) and
9           ^Longitude(longitude, event_kg:at_longitude[event])
10
11  def output = list_positions
```
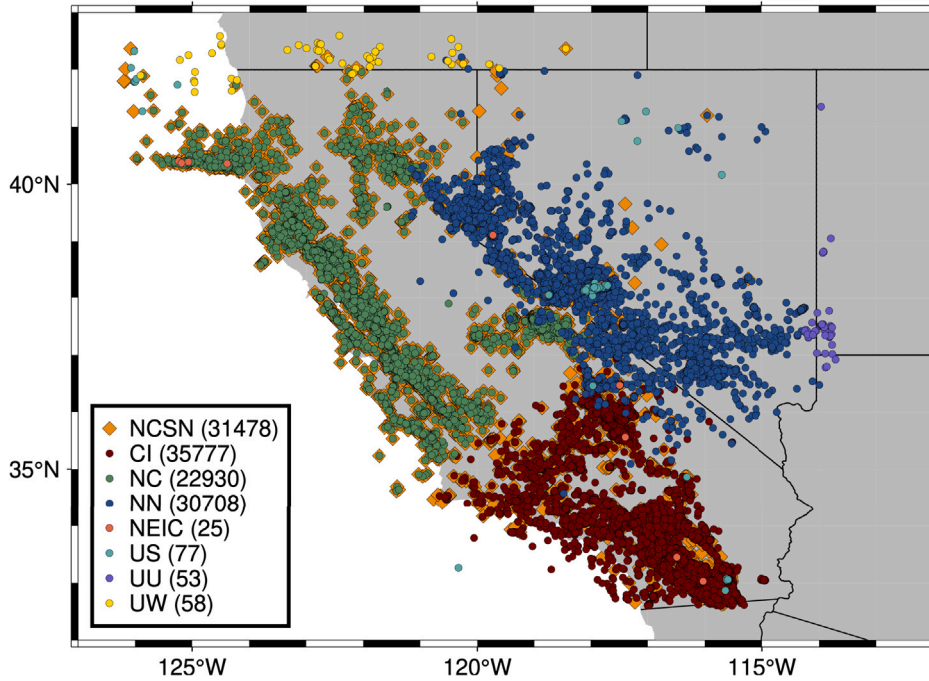


**Fig. 5.** A map of earthquake event locations found in Example 1. The epicenters of the events associated with different contributors are represented by colored symbols, as specified in the legend. The number of events from each contributor is indicated in parentheses. The original source of data for each contributor are as follows: NCSN (NCEDC), CI (USGS), NC (USGS), NN (USGS), NEIC (GCMT), US (USGS), UU (USGS), UW (USGS).

### 4.2. Example 2: An event focused study

Next, we use both event and station KGs to study a single earthquake in detail. As a case study, we examine the 2019 Ridgecrest earthquake sequence, which caused widespread shaking throughout southern California (Brandenberg et al., 2019). We aim to determine a set of strong-motion instruments that were spatially and temporally coincidental with the earthquake. The example demonstrates how to query data attributes with a KG to reduce the volume of targets for waveform data acquisition.

We use the IRIS DMC `fdsnws-station` webservice to collect station metadata around Southern California for all stations that were operational on or after the day of the earthquake. Next we obtain event data for the year 2019 from the GCMT catalog, and use both datasets to construct event and station KGs. The resulting station KG comprises ∼ 35,000 nodes and ∼ 116,000 edges in the ontology,

Listing 4: Querying the station KG (`station_kg`) through the requirements outlined in Example 2. Lines 1–2 identify the `Event Record` entity of the Ridgecrest earthquake through its event ID. Lines 5-6 traverse the graph from station to channel group to channel nodes. Line 7 filters the channels by low gain and vertically oriented instruments. Line 8 uses the relation in Listing 3 to select stations operational during the Ridgecrest earthquake. Line 9 uses relational composition to select channel groups that contain three channels, with at least one channel in the channel group satisfying Line 7. Line 10 filters event-station pairs by their epicentral distance in degrees. Relations `is_low_gain_vertical` and `event_station_radius_range` are defined in the supplementary material. From this query, latitudes and longitudes are realized similarly to Listing 2 and the query is defined in the supplementary material.

```
1   def ridgecrest_event(event) =
2       event_kg:event_id(event,^EventId["C201907060319A"])
3
4   def ridgecrest_query(station) =
5       station_kg:channel_group(station,channelgroup) and
6       station_kg:channel(channelgroup,channel) and
7       is_low_gain_vertical(channel) and
8       event_in_channel_operational_range(ridgecrest_event, channel) and
9       count(station_kg:channel[channelgroup], 3) and
10      event_station_radius_range[0.0,2.0](ridgecrest_event,station)
11      from channelgroup, channel
```

including 2316 `Station` entities. We identify the `Event Record` entity corresponding to the Mw 7.1 July 6th earthquake through its event ID, C201907060319A, obtained from the IRIS Moment Tensor page (doi:10.17611/DP/18001775) (Trabant et al., 2012).

Next, we simulate a typical query to identify scientifically useful strong-motion data relating to the 2019 event:

1. The `Station` must be within 2 degrees (222 km) of the earthquake epicenter,
2. The `Station` must contain a `Channel Group` where:

   (a) The `Channel Group` has 3 `Channel` entities,

3. The `Channel Group` must have a `Channel` where:

   (a) The `Channel` band type is either broadband or high broadband,
   (b) The `Channel` instrument type is an accelerometer,
   (c) The `Channel` is in the vertical orientation, and
   (d) The `Channel` was operational at the time of the earthquake,

Requirement (1) compares the latitude and longitude properties of the event KG (E7–8) and the stations KG (S8–9). Similarly, (3.d) compares the event KG date-time property (E6) with the station KG start and end date-time properties (S15–16), as shown in Listing 3 with query `event_in_channel_operational_range`. Requirement (2) is satisfied by summing the number of edges connecting a `Channel Group` to its `Channel` nodes. Finally, requirements (3.a–c) are accomplished by filtering the KG by `Channel` relations S11–13. We use these conditions to query for useful strong-motion stations in Listing 4, `ridgecrest_query`.

Listing 3: Querying the station KG (`station_kg`) and event KG (`event_kg`) through the constraints outlined in Example 2. This statement defines a binary relation of `Event Record` and `Channel` pairs, where the channel was operational during the event. This is accomplished through the edge relations for channel nodes `operational_from` and `operational_until` and the edge relation for event record nodes `occurred_at`. The target datetime nodes are then constrained with relational operators on Line 5.

```
1   def event_in_channel_operational_range(event, channel) =
2       station_kg:operational_from(channel, start_dt) and
3       station_kg:operational_until(channel, end_dt) and
4       event_kg:occurred_at(event, event_dt) and
5       (event_dt > start_dt) and (event_dt < end_dt)
6       from start_dt, end_dt, event_dt
```

We find that 107 stations match the query in Listing 4; their spatial distribution is shown in Fig. 6.

### 4.3. Example 3: A constrained global seismology study

To illustrate the flexibility and utility of the KG approach, we investigate a case study with highly specific constraints on station and event data: the study of the inner core through deep seismic phases (e.g., Tkalčić et al., 2013; Yu et al., 2017). Such studies commonly require high-gain, low-noise instruments with favorable orientations and often utilize earthquakes of particular magnitudes, depths, and sometimes earthquakes with high latitudes (Frost et al., 2021) or temporally repeating patterns (Yang and Song, 2023). The strictest constraint is placed upon event-station pairs, as very precise event-station epicentral distances are required to observe the necessary core-sampling seismic phases (Young et al., 2013; Tkalčić, 2015). In this example, we use KGs to efficiently determine valid event-station pairs, which informs the acquisition of waveform data for inner core studies.

This example demonstrates using one semantic model to search across a mix of semi-structured and tabular event and station data, with each schema using distinct terminology and data organizing principles. We collect station metadata from the IRIS DMC, without restrictions on geographic position, for all stations with operational channels starting on or after 2010. We use earthquake event data from the entire GCMT catalog combined with records of nuclear explosions gathered from the USGS earthquake catalog webservice. The KGs contain $\sim 1.4$ million nodes and $\sim 19$ million edges for the station KG, and $\sim 591,000$ nodes and $\sim 623,000$ edges for the event KG. In particular, we have $\sim 62,000$ `Event Record` and $\sim 49,000$ `Station` entities.

We simulate a highly constrained query to identify scientifically useful event-station pairs that could sample core phases PKPbc and PKPdf (Tkalčić, 2015). (Waveform data and its constraints, for example imposed by travel time analysis, are outside the scope of this paper.) Events are defined to have the following requirements:

1. The `Event Record` magnitude is between 5.5 and 7,
2. The `Event Record` depth is greater than 10 km, and
3. The `Event Record` latitude is either greater than 45°N or less than 45°S.

Similarly, stations have the following constraints:

4. The `Station` latitude is either greater than 45°N or less than 45°S,
5. The `Channel` band type must be either broadband or high broadband, and
6. The `Channel` instrument type must be a high-gain seismometer.

Finally, there are constraints on the event-station pairs:

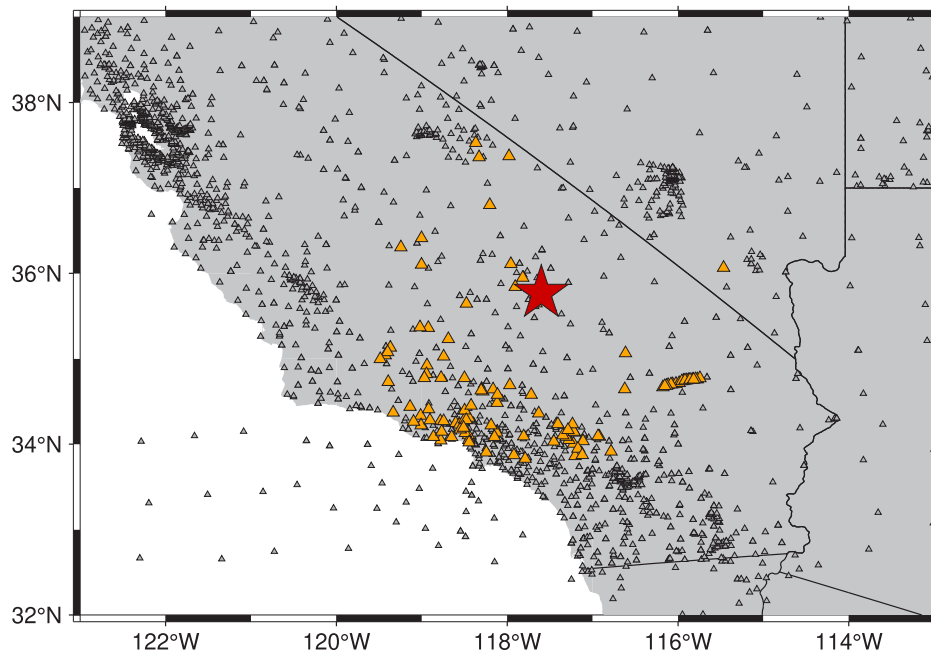7. The `Channel` must be operational at the time of the earthquake, and

**Fig. 6.** A map of stations around the 2019 Ridgecrest earthquake found in Example 2. The estimated epicenter of the earthquake is indicated with a red star (doi:10.17611/DP/18001775) (Trabant et al., 2012). The locations of all 2316 stations present in the dataset area are indicated with triangles. Stations that match the query described in Example 2 are indicated with yellow triangles. Stations that do not match the query are indicated with smaller gray triangles.

8. The epicentral distance between the Event Record epicenter and Station is in the range 147°–153°.

Requirement (7) is satisfied by the date-time comparison relation in Listing 3. Requirement (8) requires a join over at_latitude and at_longitude for both the station and event KGs, as well as calculation of the distance, and comparison with the distance range. This constraint is implemented in Listing 5 as query filter_epicentral_distance. Requirements (1–3) can be accomplished by filtering the event KG on relations derived from E5, E7, and E9 and is expressed in event_query in Listing 6. Similarly, (4–6) involve filtering the station KG on relations from S9, S11, and S12 as expressed in station_query in Listing 6. Finally, we combine these queries in the inner_core_query in Listing 6 to demand the core-sampling event-station pairs that satisfy all requirements 1–8.

Listing 5: Definition of a binary relation of Event Record and Station pairs filtered by epicentral distance, as outlined in Example 3. Lines 3-7 resolve the Event Record and Station latitudes and longitudes. The relation on Line 2 calculates an epicentral distance and is defined in the supplementary material. Line 8 compares this distance with a prescribed epicentral distance range (defined in the supplementary material). From this query, latitudes and longitudes for event-station pairs are realized similarly to Listing 2.

```
1   def filter_epicentral_distance(event,station) =
2       great_circle_distance(
3           event_kg:at_latitude[event],
4           event_kg:at_longitude[event],
5           station_kg:at_latitude[station],
6           station_kg:at_longitude[station],
7           distance
8       ) and
9       (147 < distance) and (distance < 153)
10      from distance
```

We find that 653 Event Record and 3145 Station entities match the queries event_query and station_query, respectively. Of the possible ∼ 2 million event-station pairs, only ∼ 127,000 match with inner_core_query. The spatial distributions of a sample of events, stations, and paths are shown in Fig. 7, revealing areas for potential inner core studies.

## 5. Discussion and conclusions

In this paper, we introduce KGs for semantic modeling of seismic station and event data. We define ontologies reflecting domain knowledge in seismology and present three examples of how knowledge from schema-diverse, real-world data can be used to construct KGs. Our examples illustrate how KGs de-emphasize schema-related details of the data, allowing a focus on composition of intelligible queries for data exploration and analysis.

We see several promising avenues for future applications of KGs in seismology. A natural progression would investigate the representation of seismic waveform data, which could be represented in a relational KG using hypergraphs. KGs could be particularly applicable to dense, highly relational, temporary deployments, such as digital acoustic seismometry (Lindsey and Martin, 2021), ocean-bottom seismometry (Suetsugu and Shiobara, 2014), or controlled source seismometry (Mondol, 2010). Another natural step would test the construction of KGs from other seismic data-formats, such as (dataless) SEED (Ahern et al., 2009), or QuakeML (Schorlemmer et al., 2011). The KGs presented here can be expanded to incorporate knowledge that seismologists may wish to model. Potential extensions could add properties for event focal mechanisms, centroid parameters, or include entities for instrument response (Ringler and Bastien, 2020) or virtual networks (Ahern, 2004). For instance, earthquake parameters reported by catalogs are occasionally revised following updates to moment tensor inversions (Weatherill et al., 2016); this could be modeled in our ontology by including additional attributes to uniquely identify non-preferred event record nodes, and adding edges connecting them to the preferred event record.

Alongside technical developments, we see potential for integrating seismic KGs with other geoscience products. Combination with ontology-driven knowledge models of geological maps—e.g., Mantovani et al. (2020)—could enable reasoning concerning the local geology around a station. A KG approach may be amenable for high-level, data-rich seismic products, including ShakeMaps (Worden et al., 2010), "Did You Feel It?" maps (Wald et al., 2011), seismic velocity models (Ritsema and Lekić, 2020), or even Green's function databases (van Driel et al., 2015). Finally, there is potential for integrating KGs into modern seismic machine learning methodologies (e.g., Zhu and Beroza, 2019;

Listing 6: Querying the event and station KGs through the constraints outlined in Example 3. The relation `event_query` defines a unary relation of `Event Record` entities that match the given constraints: Line 2 matches with events below a specified depth; Line 3 matches with events in a range of magnitudes; and Lines 4–5 similarly match with events in two disjoint ranges of latitudes. The relation `station_query` defines a similar unary relation of `Station` entities that match the given constraints: Lines 8–9 traverses the graph from `Station` to `Channel`; Lines 10–12 filter the `Channel` for (high-)broadband and high-gain seismometers; and Lines 13–14 match with stations in two disjoint ranges of latitudes. The relation `inner_core_query` combines the two previous relations in Lines 18–19, and further constrains with `event_in_channel_operational_range` (Line 20) and `filter_epicentral_distance` (Line 21), (previously defined in Listings [3] and [5], respectively. Auxiliary relations used here are defined in the supplementary material.

```
1   def event_query(event) =
2       depth_below(10,event) and
3       filter_event_magnitude(5.5,7,event) and
4       ( filter_latitude(event_kg,55,90,event) or
5         filter_latitude(event_kg,-90,-55,event) )
6
7   def station_query(station, channel) =
8       station_kg:channel_group(station,channelgroup) and
9       station_kg:channel(channelgroup,channel) and
10      ( station_kg:band_type(channel,^BandType["B"]) or
11        station_kg:band_type(channel,^BandType["H"]) ) and
12      station_kg:instrument_type(channel,^InstrumentType["H"]) and
13      ( filter_latitude(station_kg,55,90,station) or
14        filter_latitude(station_kg,-90,-55,station) )
15      from channelgroup
16
17  def inner_core_query(event,station) =
18      event_query(event) and
19      station_query(station,channel) and
20      event_in_channel_operational_range(event,channel) and
21      filter_epicentral_distance(event,station)
22      from channel
```
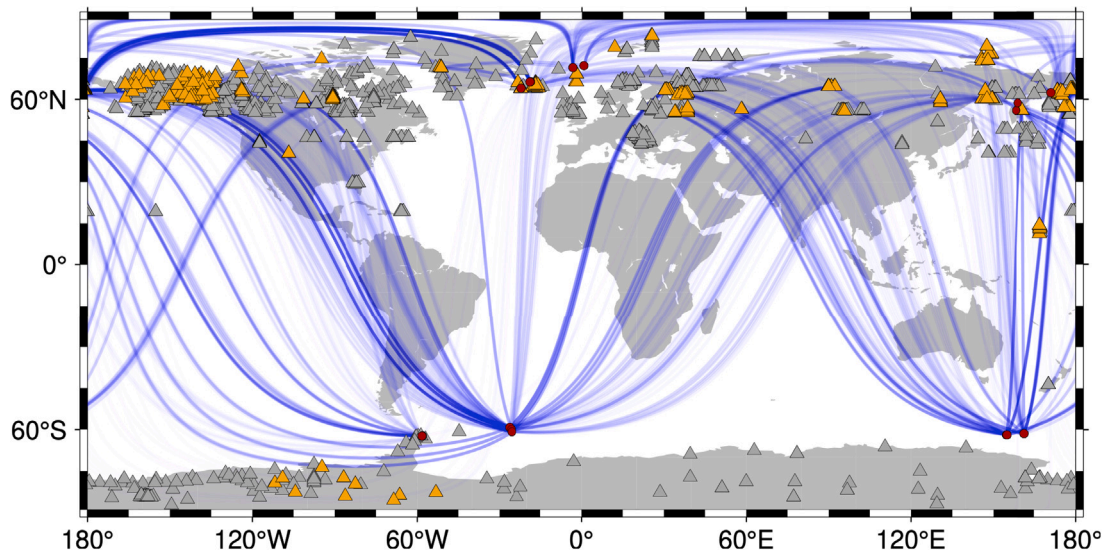


**Fig. 7.** A map of events and stations found in Example 3. For visual clarity, we restrict this plot to only show events occurring in 2020: 16 of the total 653, corresponding to 3840 event-station pairs. The locations of earthquakes and nuclear explosions satisfying `event_query` are indicated with circles. The locations of stations satisfying `station_query` are indicated with triangles. Gray symbols indicate events or stations which match `event_query` or `station_query` but do not appear in any of the valid event-station pairs defined by `inner_core_query`. Colored symbols indicate events and stations which satisfy `inner_core_query`, forming valid event-station pairs. Great-circle paths between the valid event-station pairs are shown as transparent blue lines.

Yeck et al., 2021). In addition to providing a semantic layer to aid explainability (Lecue, 2020), KGs could enhance the training of seismic machine learning models by embedding prior logic and structure into training data (Hogan et al., 2021).

In conclusion, we believe that KGs have a promising interconnected and interdisciplinary future in seismology. Used as complementary tools to augment traditional seismic databases, KGs offer flexibility and accessibility. In this application, they are best utilized when provided by institutional data providers or large research groups, rather than individual researchers. We look forward to further exploring the potential of KGs in seismology and beyond.

**Data, resources, and code availability**

The data underlying this paper are available in the Dryad Digital Repository, at https://doi.org/10.6078/D1P430, and the Zenodo open data repository, at https://doi.org/10.5281/zenodo.8346843. All code used in this paper is available in the Zenodo open data repository, at https://doi.org/10.5281/zenodo.10183012.

The resources mentioned in the article and their corresponding references: International Federation of Digital Seismograph Networks (FDSN) Station Extended Markup Language format/specification (StationXML) is available at https://www.fdsn.org/xml/station/; FDSN

Source Identifiers specification is available at http://docs.fdsn.org/projects/source-identifiers/.

Data sources used in this article and their corresponding references: station metadata from the Incorporated Research Institutions for Seismology (IRIS), acquired using the `fdsnws-station` webservice at https://service.iris.edu/fdsnws/station/1/; earthquake event data from the Global Centroid-Moment Tensor (GCMT) catalog, acquired using the webservice at https://www.globalcmt.org; earthquake event data from the Northern California Earthquake Data Center (NCEDC), acquired using the Northern California Earthquake Catalog Search webservice at https://doi.org/10.7932/NCEDC; and earthquake event data from the United States Geological Survey (USGS) Advanced National Seismic System (ANSS) Comprehensive Earthquake Catalog (ComCat), acquired using the webservice at https://doi.org/10.5066/F7MS3QZH. Figs. 3 and 4 were created with `ormjs`, available at https://github.com/crhunt/ormjs. All websites were last accessed in August 2023.

## CRediT authorship contribution statement

**William Davis:** Writing – original draft, Visualization, Validation, Software, Investigation, Data curation, Conceptualization. **Cassandra R. Hunt:** Writing – review & editing, Supervision, Software.

## Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Data availability

The data and code associated with this paper are available from the DataDryad and Zenodo links found in the Data, Resources, and Code Availability section.

## Acknowledgments

## References

Abu-Salih, B., 2021. Domain-specific knowledge graphs: A survey. J. Netw. Comput. Appl. 185, 103076.

Ahern, T., 2004. Virtual networks defined. URL: http://ds.iris.edu/ds/newsletter/vol6/no2/.

Ahern, T., Casey, R., Barnes, D., Benson, R., Knight, T., Trabant, C., 2009. SEED reference manual, version 2.4. URL: http://www.fdsn.org/seed_manual/SEEDManual_V2.4.pdf.

Arrais, S., Urquiza-Aguiar, L., Tripp-Barba, C., 2022. Analysis of information availability for seismic and volcanic monitoring systems: A review. Sensors 22 (14), 5186.

Arrowsmith, S.J., Trugman, D.T., MacCarthy, J., Bergen, K.J., Lumley, D., Magnani, M.B., 2022. Big data seismology. Rev. Geophys. 60 (2), e2021RG000769.

Benson, R.B., Ronan, T., Suleiman, Y.Y., Casey, R.E., Trabant, C.M., Templeton, M., Carter, J., 2019. An Introduction to the StationXML-SEED-converter and StationXML-validator, a Set of FDSN-StationXML Metadata Utilities. In: AGU Fall Meeting 2019. AGU, NS21B–0816.

Beyreuther, M., Barsch, R., Krischer, L., Megies, T., Behr, Y., Wassermann, J., 2010. ObsPy: A Python toolbox for seismology. Seismol. Res. Lett. 81 (3), 530–533.

Brandenberg, S.J., Wang, P., Nweke, C.C., Hudson, K., Mazzoni, S., Bozorgnia, Y., Hudnut, K.W., Davis, C.A., Ahdi, S.K., Zareian, F., et al., 2019. Preliminary Report on Engineering and Geological Effects of the July 2019 Ridgecrest Earthquake Sequence. Technical Report, Geotechnical Extreme Event Reconnaissance Association.

Buland, R., 2012. Seismic station codes – new coding standards. In: Bormann, P. (Ed.), New Manual of Seismological Observatory Practice 2. NMSOP-2, Deutsches GeoForschungsZentrum (GFZ), Potsdam, pp. 1–9. http://dx.doi.org/10.2312/GFZ.NMSOP-2_IS_10.3.

Date, C.J., 2006. The Relational Database Dictionary: A Comprehensive Glossary of Relational Terms and Concepts, with Illustrative Examples. O'Reilly Media, Inc.

Dost, B., Zednik, J., Havskov, J., Willemann, R., Bormann, P., 2009. Seismic data formats, archival and exchange. In: New Manual of Seismological Observatory Practice. NMSOP, Deutsches GeoForschungsZentrum GFZ, pp. 1–20.

Dziewonski, A.M., Chou, T.-A., Woodhouse, J.H., 1981. Determination of earthquake source parameters from waveform data for studies of global and regional seismicity. J. Geophys. Res.: Solid Earth 86 (B4), 2825–2852.

Ekström, G., Nettles, M., Dziewoński, A., 2012. The global CMT project 2004–2010: Centroid-moment tensors for 13,017 earthquakes. Phys. Earth Planet. Inter. 200, 1–9.

Falco, N., Clark, A., Trabant, C., 2017. WILBER and PyWEED: Event-based seismic data request tools. In: AGU Fall Meeting Abstracts. 2017. pp. T44D–06.

Frost, D.A., Lasbleis, M., Chandler, B., Romanowicz, B., 2021. Dynamic history of the inner core constrained by seismic anisotropy. Nat. Geosci. 14 (7), 531–535.

García-González, H., Boneva, I., Staworko, S., Labra-Gayo, J.E., Lovelle, J.M.C., 2020. ShExML: improving the usability of heterogeneous data mapping languages for first-time users. PeerJ Comput. Sci. 6, e318.

Gil, Y., Pierce, S.A., Babaie, H., Banerjee, A., Borne, K., Bust, G., Cheatham, M., Ebert-Uphoff, I., Gomes, C., Hill, M., et al., 2018. Intelligent systems for geosciences: an essential research agenda. Commun. ACM 62 (1), 76–84.

Gutiérrez, C., Sequeda, J.F., 2021. Knowledge graphs. Commun. ACM 64 (3), 96–104.

Halevy, A., Rajaraman, A., Ordille, J., 2006. Data integration: The teenage years. In: Proceedings of the 32nd International Conference on Very Large Data Bases. pp. 9–16.

Halpin, T., 2015. Object-Role Modeling Fundamentals: A Practical Guide to Data Modeling with ORM. Technics Publications.

Hofer, M., Obraczka, D., Saeedi, A., Köpcke, H., Rahm, E., 2023. Construction of knowledge graphs: State and challenges. arXiv preprint arXiv:2302.11509.

Hogan, A., Blomqvist, E., Cochez, M., d'Amato, C., Melo, G.d., Gutierrez, C., Kirrane, S., Gayo, J.E.L., Navigli, R., Neumaier, S., et al., 2021. Knowledge graphs. ACM Comput. Surv. 54 (4), 1–37.

Hölsch, J., Schmidt, T., Grossniklaus, M., 2017. On the performance of analytical and pattern matching graph queries in neo4j and a relational database. In: EDBT/ICDT 2017 Joint Conference: 6th International Workshop on Querying Graph Structured Data. GraphQ.

Humphries, B., 2021. Relational paradigm. URL: https://relational.ai/blog/relational-paradigm.

Hutko, A.R., Bahavar, M., Trabant, C., Weekly, R.T., Fossen, M.V., Ahern, T., 2017. Data products at the IRIS-DMC: Growth and usage. Seismol. Res. Lett. 88 (3), 892–903.

International Seismological Centre (ISC), 2020. International Seismograph Station Registry (IR). International Seismological Centre Thatcham, http://dx.doi.org/10.31905/EL3FQQ40.

Krischer, L., Smith, J., Lei, W., Lefebvre, M., Ruan, Y., de Andrade, E.S., Podhorszki, N., Bozdağ, E., Tromp, J., 2016. An adaptable seismic data format. Geophys. J. Int. 207 (2), 1003–1011.

Lecue, F., 2020. On the role of knowledge graphs in explainable AI. Semantic Web 11 (1), 41–51.

Lindsey, N.J., Martin, E.R., 2021. Fiber-optic seismology. Ann. Rev. Earth Planet. Sci. 49, 309–336.

Ma, X., 2022. Knowledge graph construction and application in geosciences: A review. Comput. Geosci. 105082.

Ma, X., Fox, P., Rozell, E., West, P., Zednik, S., 2014. Ontology dynamics in a data life cycle: Challenges and recommendations from a Geoscience Perspective. J. Earth Sci. 25, 407–412.

Mantovani, A., Piana, F., Lombardo, V., 2020. Ontology-driven representation of knowledge for geological maps. Comput. Geosci. 139, 104446.

Mohammadpoor, M., Torabi, F., 2020. Big Data analytics in oil and gas industry: An emerging trend. Petroleum 6 (4), 321–328.

Mondol, N.H., 2010. Seismic exploration. Petroleum Geosci. 1, 375–402.

Monteiro, J., Sá, F., Bernardino, J., 2023. Experimental evaluation of graph databases: JanusGraph, nebula graph, Neo4j, and TigerGraph. Appl. Sci. 13 (9), 5770.

NCEDC, 2014. Northern California earthquake data center, UC Berkeley seismological laboratory. http://dx.doi.org/10.7932/NCEDC.

Newman, R., Clark, A., Trabant, C., Karstens, R., Hutko, A., Casey, R., Ahern, T., 2013. Wilber 3: A Python-Django web application for acquiring large-scale event-oriented seismic data. In: Agu Fall Meeting Abstracts, 2013. IN51B–1543.

Obraczka, D., Schuchart, J., Rahm, E., 2021. EAGER: Embedding-assisted entity resolution for knowledge graphs. arXiv preprint arXiv:2101.06126.

RelationalAI (RAI), 2021a. The rel language. URL: https://docs.relational.ai/rel.

RelationalAI (RAI), 2021b. The relational knowledge graph system (RKGS). URL: https://docs.relational.ai/rkgms.

Ringler, A.T., Anthony, R.E., Aster, R., Ammon, C., Arrowsmith, S., Benz, H., Ebeling, C., Frassetto, A., Kim, W.-Y., Koelemeijer, P., et al., 2022. Achievements and prospects of global broadband seismographic networks after 30 years of continuous geophysical observations. Rev. Geophys. (1985) 60 (3).

Ringler, A.T., Bastien, P., 2020. A brief introduction to seismic instrumentation: Where does my data come from? Seismol. Res. Lett. 91 (2A), 1074–1083.

Ritsema, J., Lekić, V., 2020. Heterogeneity of seismic wave velocity in Earth's mantle. Ann. Rev. Earth Planet. Sci. 48, 377–401.

Schorlemmer, D., Euchner, F., Kästli, P., Saul, J., Group, Q.W., et al., 2011. QuakeML: status of the XML-based seismological data exchange format. Ann. Geophys. 54 (1).

Southern California Earthquake Data Center (SCEDC), 2021. The SCEDC Earthquake Data AWS Public Dataset, Bucket s3://scedc-pds; us-west-2. URL: https://scedc.caltech.edu/data/cloud.html.

Spica, Z.J., Ajo-Franklin, J., Beroza, G.C., Biondi, B., Cheng, F., Gaite, B., Luo, B., Martin, E., Shen, J., Thurber, C., et al., 2023. PubDAS: A public distributed acoustic sensing datasets repository for geosciences. Seismol. Soc. Am. 94 (2A), 983–998.

Stirewalt, R., 2022. Experience report: building enterprise applications using LogiQL and Rel. URL: https://www.hytradboi.com/2022/experience-report-building-enterprise-applications-using-logiql-and-rel.

Stirewalt, R., Búr, M., 2023. The RAI way: A technical analysis and design method for building enterprise semantic layers. In: International Conference on Advanced Information Systems Engineering. Springer, pp. 74–79.

Suetsugu, D., Shiobara, H., 2014. Broadband ocean-bottom seismology. Ann. Rev. Earth Planet. Sci. 42, 27–43.

Sun, Z., Zhang, Q., Hu, W., Wang, C., Chen, M., Akrami, F., Li, C., 2020. A benchmarking study of embedding-based entity alignment for knowledge graphs. arXiv preprint arXiv:2003.07743.

Timón-Reina, S., Rincón, M., Martínez-Tomás, R., 2021. An overview of graph databases and their applications in the biomedical domain. Database 2021, baab026.

Tkalčić, H., 2015. Complex inner core of the Earth: The last frontier of global seismology. Rev. Geophys. 53 (1), 59–94.

Tkalčić, H., Young, M., Bodin, T., Ngo, S., Sambridge, M., 2013. The shuffling rotation of the Earth's inner core revealed by earthquake doublets. Nat. Geosci. 6 (6), 497–502.

Trabant, C., Benson, R.B., Carter, J., Casey, R.E., 2019. The evolution of seismological data standards and what the changes mean for users. In: AGU Fall Meeting Abstracts, 2019. pp. S21H–0632.

Trabant, C., Hutko, A.R., Bahavar, M., Karstens, R., Ahern, T., Aster, R., 2012. Data products at the IRIS DMC: Stepping stones for research and other applications. Seismol. Res. Lett. 83 (5), 846–854.

Trugman, D.T., Fang, L., Ajo-Franklin, J., Nayak, A., Li, Z., 2022. Preface to the focus section on big data problems in seismology. Seismol. Soc.. Am. 93 (5), 2423–2425.

United States Geological Survey (USGS), 2021. US Geological Survey 21st-century science strategy 2020–2030. http://dx.doi.org/10.3133/cir1476.

United States Geological Survey (USGS), Earthquake Hazards Program, 2017. Advanced National Seismic System (ANSS) comprehensive catalog of earthquake events and products: Various. http://dx.doi.org/10.5066/F7MS3QZH.

van Driel, M., Krischer, L., Stähler, S.C., Hosseini, K., Nissen-Meyer, T., 2015. Instaseis: Instant global seismograms based on a broadband waveform database. Solid Earth 6 (2), 701–717.

Wald, D.J., Quitoriano, V., Worden, C.B., Hopper, M., Dewey, J.W., 2011. USGS "Did You Feel It?" internet-based macroseismic intensity maps. Ann. Geophys. 54 (6), 688.

Wang, R.Y., Strong, D.M., 1996. Beyond accuracy: What data quality means to data consumers. J. Manag. Inform. Syst. 12 (4), 5–33.

Weatherill, G., Pagani, M., Garcia, J., 2016. Exploring earthquake databases for the creation of magnitude-homogeneous catalogues: tools for application on a regional and global scale. Geophys. J. Int. 206 (3), 1652–1676.

Weertman, B., 2010. Web services at the DMC, *IRIS Data Services Newsletter*. URL: https://ds.iris.edu/ds/newsletter/vol12/no3/44/web-services-at-the-dmc/.

Wing, J.M., 2019. The data life cycle. Harvard Data Sci. Rev. 1 (1), 6.

Worden, C., Wald, D., Allen, T., Lin, K., Garcia, D., Cua, G., 2010. A revised ground-motion and intensity interpolation scheme for ShakeMap. Bull. Seismol. Soc. Am. 100 (6), 3083–3096.

Xiao, G., Ding, L., Cogrel, B., Calvanese, D., 2019. Virtual knowledge graphs: An overview of systems and use cases. Data Intell. 1 (3), 201–223.

Yang, Y., Song, X., 2023. Multidecadal variation of the Earth's inner-core rotation. Nat. Geosci. 1–6.

Yeck, W.L., Patton, J.M., Ross, Z.E., Hayes, G.P., Guy, M.R., Ambruz, N.B., Shelly, D.R., Benz, H.M., Earle, P.S., 2021. Leveraging deep learning in global 24/7 real-time earthquake monitoring at the National Earthquake Information Center. Seismol. Soc. Am. 92 (1), 469–480.

Young, M., Tkalčić, H., Bodin, T., Sambridge, M., 2013. Global P wave tomography of Earth's lowermost mantle from partition modeling. J. Geophys. Res.: Solid Earth 118 (10), 5467–5486.

Yu, E., Bhaskaran, A., Chen, S.-L., Ross, Z.E., Hauksson, E., Clayton, R.W., 2021. Southern California earthquake data now available in the AWS Cloud. Seismol. Res. Lett. 92 (5), 3238–3247.

Yu, W.-c., Su, J., Song, T.-R.A., Huang, H.-H., Mozziconacci, L., Huang, B.-S., 2017. The inner core hemispheric boundary near 180° W. Phys. Earth Planet. Inter. 272, 1–16.

Zhan, Z., 2020. Distributed acoustic sensing turns fiber-optic cables into sensitive seismic antennas. Seismol. Res. Lett. 91 (1), 1–15.

Zhu, W., Beroza, G.C., 2019. PhaseNet: A deep-neural-network-based seismic arrival-time picking method. Geophys. J. Int. 216 (1), 261–273.