

1 Artificial intelligence for ocean science data integration:
2 current state, gaps, and way forward

3 Tomer Sagi^{1,2,*}, Yoav Lehahn^{3,*}, Koby Bar¹

4 ¹Department of Information Systems, University of Haifa, Haifa, Israel

5 ²Department of Computer Science, Aalborg University, Aalborg, Denmark.

6 ³Department of Marine Geosciences, Charney School of Marine Sciences, University of Haifa,
7 Haifa, Israel.

8 *tsagi@is.haifa.ac.il,ylehahn@univ.haifa.ac.il

9 **Abstract**

10 Oceanographic research is a multidisciplinary endeavor that involves the acqui-
11 sition of an increasing amount of in-situ and remotely sensed data. A large and
12 growing number of studies and data repositories are now available on-line. How-
13 ever, manually integrating different datasets is a tedious and grueling process lead-
14 ing to a rising need for automated integration tools. A key challenge in oceano-
15 graphic data integration is to map between data sources that have no common
16 schema and that were collected, processed, and analyzed using different method-
17 ologies. Concurrently, artificial agents are becoming increasingly adept at extract-
18 ing knowledge from text and using domain ontologies to integrate and align data.
19 Here, we deconstruct the process of ocean science data integration, providing a
20 detailed description of its three phases: discover, merge, and evaluate/correct. In
21 addition, we identify the key missing tools and underutilized information sources
22 currently limiting the automation of the integration process. The efforts to address
23 these limitations should focus on (i) development of artificial intelligence-based
24 tools for assisting ocean scientists in aligning their schema with existing ontolo-
25 gies when organizing their measurements in datasets; (ii) extension and refinement
26 of conceptual coverage of – and conceptual alignment between – existing ontolo-
27 gies, to better fit the diverse and multidisciplinary nature of ocean science; (iii)
28 creation of ocean-science-specific *entity resolution* benchmarks to accelerate the
29 development of tools utilizing ocean science terminology and nomenclature; (iv)
30 creation of ocean-science-specific schema matching and mapping benchmarks to
31 accelerate the development of matching and mapping tools utilizing semantics en-

32 coded in existing vocabularies and ontologies; (v) annotation of datasets, and de-
33 velopment of tools and benchmarks for the extraction and categorization of data
34 quality and preprocessing descriptions from scientific text; and (vi) creation of
35 large-scale word embeddings trained upon ocean science literature to accelerate
36 the development of information extraction and matching tools based on artificial
37 intelligence.

38 **1. Introduction**

39 The study of the ocean is one of the biggest scientific challenges of the 21st cen-
40 tury. It has a direct impact on our understanding of Earth’s climate (Stocker et al.,
41 2013) and biogeochemical cycling (Field et al., 1998), as well as on our ability to
42 provide human society with food, chemicals, and energy (Lehahn et al., 2016).
43 Oceanographic research relies mainly on in-situ and remotely-sensed observa-
44 tions, which describe physical, chemical, and biological seawater properties at
45 a given time and place. These observations are collected from various crewed and
46 autonomous platforms, including research vessels, floats (Roemmich et al., 2009),
47 drifters (Lumpkin et al., 2017), autonomous vehicles (Eriksen et al., 2001), and
48 satellites (Lehahn et al., 2018), providing an abundance of interdisciplinary infor-
49 mation on processes occurring over a wide range of spatial (from micrometers to
50 thousands of kilometers) and temporal (from seconds to decades) scales.

51 Over the last century, numerous in-situ and remotely-sensed measurements
52 have been taken, resulting in the creation of an increasingly large amount of
53 oceanic data. In recent years, with the enhanced utilization of satellites and au-
54 tonomous observation platforms, these data are collected at a blistering rate. Im-
55 proving the scientific community’s ability to integrate, share, and explore this vast
56 amount of data is an urgent task that will contribute substantially to our under-
57 standing of the ocean and its role in the Earth system.

58 Several public data repositories have emerged to enable the archiving and
59 sharing of data collected between researchers. For example, PANGEA (2020),
60 a data repository for publishing and distributing georeferenced data from Earth
61 system research, hosts more than 370,000 datasets. The National Centers for En-
62 vironmental Information (NCEI, 2020) stores over 25 petabytes of atmospheric,
63 coastal, oceanic, and geophysical data. Copernicus (European Commission, 2020)
64 archives datasets from several domains such as marine, climate, and agriculture,
65 as part of a European Union program for observing the Earth. The extensive avail-
66 ability of data repositories provides oceanographic researchers with the ability to
67 tap into a multitude of data collected by their peers and use it in their own studies.

68 One of the main obstacles for a researcher when compiling data from existing data
69 sources is to overcome the semantic distance between datasets. Thus, when con-
70 ducting such research, there is a need for manual data integration work done by
71 an expert. In a recent review (Gregory et al., 2019), the authors described some of
72 the challenges facing researchers when manually integrating data from multiple
73 disparate studies.

74 Data integration is the art and science of reconciling two or more collections
75 of data with each other. Data integration is as old as data. In 1975, the National
76 Bureau of Standards and the Association for Computing Machinery issued the rec-
77 ommendation that, when integrating data from digital and physical files into the
78 newly standardized Database Management Systems (DBMS), practitioners should
79 maintain a data-dictionary to enable efficient and effective data integration (Berg,
80 1976). With the emergence of the federated database (Hammer and McLeod,
81 1979), a database composed of multiple independent database systems, the need
82 for a central mediated schema was created. A secondary problem created by fed-
83 erated databases was the prevalence of unwanted data duplication between the
84 systems. The advent, and subsequent popularity of the World Wide Web, brought
85 about a host of new opportunities for sharing data, providing portals and services
86 based on the integration of data from multiple sources covering the same do-
87 main, such as the domain of travel reservation (e.g., www.orbitz.com). The
88 process of data integration began as a manual one (Goodhue et al., 1992), gradu-
89 ally transitioning to a semi-automated process supported by software tools. The
90 arrival of Big Data has increased both the number and sizes of available data-
91 sources, bringing about additional challenges and opportunities for data integra-
92 tion (Dong and Srivastava, 2015).

93 We are at a time where artificial intelligence (AI) is applied ubiquitously across
94 scientific domains and disciplines. First and foremost of AI research fields is the
95 field of machine learning (ML), the science of building software that learns from
96 experience. Recent years have seen a concurrent increase in data (serving as expe-
97 rience for ML) and available cloud computing solutions to utilize the data. These
98 phenomena, together with the arrival of deep learning (DL) as an efficient and
99 effective method for ML, have caused ML to expand into an increasing number
100 of fields (Jordan and Mitchell, 2015). Pioneered by Doan et al. (2002), the use of
101 ML in data integration has been expected for some time now (Halevy et al., 2006).
102 Recently, widespread use of ML in data integration appears to be the new norm
103 (see review by Dong and Rekatsinas, 2018). Concretely, ML has been used to cre-
104 ate weighted ensembles of schema matchers (Gal and Sagi, 2010), map relational
105 databases into ontologies (De Uña et al., 2018), and create sub-groups of records

106 to speed up entity resolution (see review by O’Hare et al., 2019). However, the ad-
107 vances in data integration and specifically AI-assisted data integration have been
108 utilized sparingly, if at all, in the ocean sciences.

109 In this paper, we systematically deconstruct the process of integrating a mul-
110 titude of datasets in the ocean science domain into specific phases and tasks. For
111 each task, we review state of the art in AI-assisted data integration and discuss
112 the barriers and challenges to its adoption in the ocean sciences. We begin in the
113 following section by formally defining and providing background on artificial in-
114 telligence, data integration, and how they are used together. We then present our
115 model of data integration processes in ocean science and how artificial intelli-
116 gence can support these efforts. To demonstrate the implications of having ocean-
117 science-specific-AI tools, we then describe and provide results from an automated
118 entity extraction task on oceanic datasets.

119 **2. Background and definitions**

120 Before we dive into the use of artificial intelligence for data integration in ocean
121 sciences, we review data integration (DI), artificial intelligence (AI), and the use
122 of AI techniques in DI.

123 *2.1. Data integration*

124 DI is the process of combining two or more datasets. Datasets are collections of
125 structured data described by a *data description*, also known as a *schema*. A dataset
126 may be simple as a table, with rows as data and the header row as a schema, or
127 complex as a netCDF (UNIDATA, 2019) file containing numerical matrices.

128 Figure 1 reviews the five components of the DI process. *Schema match-*
129 *ing* (1) aligns two or more schemas to find correspondences between them (see
130 survey by Shvaiko and Euzenat, 2013 and books: Gal, 2011; Bellahsene et al.,
131 2011). *Schema mapping* (2) operationalizes these correspondences into data-
132 transformation functions (e.g., Alexe et al., 2011). (3) *Entity resolution* is the
133 task of identifying different instances related to the same entity (see surveys
134 Papadakis et al., 2016; O’Hare et al., 2019). *Entity consolidation* (4) is the process
135 of merging all data about the same entity coherently (e.g., Hogan et al., 2012). An
136 orthogonal but crucial component of the DI process is *data cleansing* (5), which
137 can be applied to both the original data and the merged dataset (Abedjan et al.,
138 2016).

139 Note that entity consolidation is designed for database records, where each
140 property has a single value. Most oceanic datasets are comprised of both database-

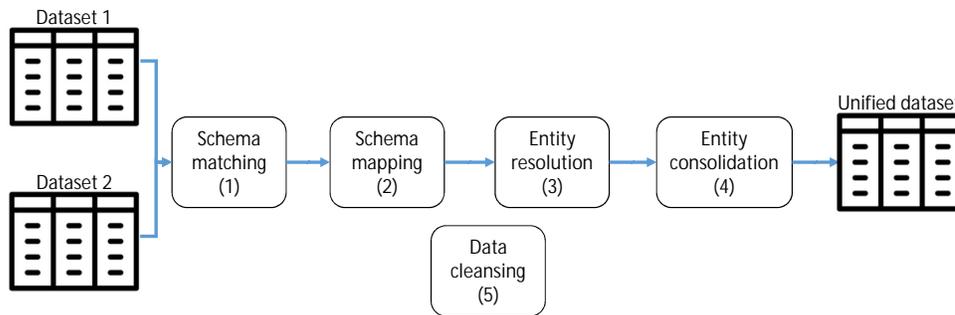


Figure 1. The process of data integration.

The data integration process takes two datasets and combines them into a unified dataset by performing five composable tasks. Schema matching (1) aligns the schemas of the two datasets. Schema mapping (2) performs any transformations required by the different semantic of the aligned fields. Entity resolution (3) identifies duplicate records, and entity consolidation (4) merges them. Data cleansing (5) can be applied at any point to detect and correct errors.

141 style records recording a dataset’s metadata and a large series of numbers varying
 142 over geographical or temporal dimensions. Integrating the numerical component,
 143 introduces two new dimensions to the integration process, namely resolution and
 144 distance. Numerical analysis and model building requires a continuous set of num-
 145 bers with the same resolution. For example, satellite images might have a spatial
 146 resolution of 1 km and a temporal resolution of one day, while a buoy in the same
 147 area and time has a pinpoint spatial resolution but may often lay a few kilometers
 148 away from the nearest sea surface image edge, due to cloud cover. To build an
 149 integrated model over both sets, one must perform interpolation and extrapolation
 150 and assess the reliability of their model given these differences and the methods
 151 employed to bridge them. Multi-sensor data fusion techniques (Lesiv et al., 2016;
 152 Waltz and Waltz, 2017) have diversified and grown from statistically based meth-
 153 ods to more elaborate ML-based methods. In the interest of brevity and focus, we
 154 limit the exploration of this task in the rest of this paper, leaving it for future work.

155 **Example 1 Schema matching and mapping (Figure. 2).** A researcher wishes to
 156 integrate PANGAEA dataset 759517 (Semina and Mikaelyan, 1994) with dataset
 157 2690 stored on EDMED (British Oceanographic Data Centre, 2020). Figure 2
 158 presents the correspondences between the two datasets’ schemas, a result of a
 159 manual schema matching process. A note added to the Nitrate field of the PAN-

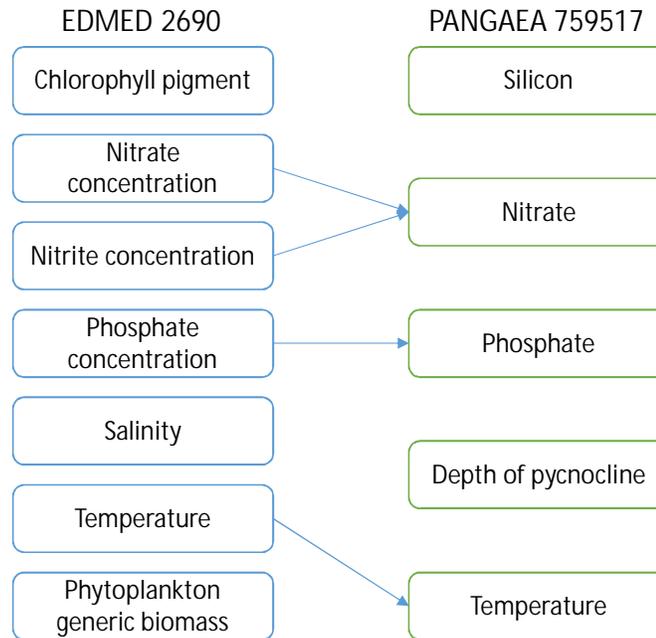


Figure 2. Schema matching of two oceanic datasets.

The figure shows correspondences created by a schema matching process between the schema of an EDMED dataset and that of a PANGAEA dataset.

160 *GAEA dataset identifies this field as actually measuring the sum of nitrates and*
 161 *nitrites, justifying the correspondences to the Nitrite and Nitrate fields in the*
 162 *EDMED dataset. This double correspondence can be converted later to a sum*
 163 *of the two values in the schema mapping process to convert data points from these*
 164 *fields under the EDMED schema to the PANGAEA schema.*

165 **Example 2 Entity resolution.** *Consider Table 1 where the same data point is pre-*
 166 *sented from the diatom data integration effort by Leblanc et al. (2012) (first row)*
 167 *and one of its constituent datasets, a supplement to Assmy et al. (2007) (second*
 168 *row). We manually schema-matched and mapped the second row to the first row's*
 169 *schema; however, it is still unclear if indeed, these represent the same data point.*
 170 *For large datasets, the entity resolution task may be daunting, requiring n^2 com-*
 171 *parisons where n is the number of records over all datasets. Thus, common ap-*
 172 *proaches perform a process of blocking, where records are grouped by (one or*
 173 *more) shared properties. In our example, these two data points were part of a*

Table 1. Entity resolution: two records mapped into the same schema

Project ID	Cruise or station ID	Date	Longitude	Latitude	Name entry
EISENEX	out of +Fe patch st° 108	11-29-2000	20.60	-47.67	<i>Thalassionema nitzschoides</i> < 20 μ m
European iron enrichment experiment in the Southern Ocean (EisenEx)	PS58/108-1 (CTD149)	2000-11-29T15:33:00	20.64733	-47.66817	<i>Thalassionema nitzschioides</i> var. <i>lanceolata</i> , biomass as carbon [μ g/l] (<i>T. nitzschioides</i> var. <i>lanceolata</i> C)

174 dataset containing 293,000 data points, of which more than half may be dupli-
 175 cates. To avoid performing 8.6×10^{11} comparisons, we could first group records
 176 by the longitude, latitude, depth, and date, and then perform comparisons only
 177 within each group (block in entity resolution terms).

178 Entity resolution can occur at different levels of granularity and for different
 179 entities appearing in the dataset. The example given above identified the same
 180 data item in the two datasets, similarly, the authors were required to resolve dif-
 181 ferent diatom species described differently. In the authors’ own words: “In total,
 182 1364 different taxonomic entries were found, but were reduced to 727 different
 183 taxonomic lines...”

184 2.2. Artificial intelligence

185 Kaplan and Haenlein (2019) define AI as: “a system’s ability to interpret exter-
 186 nal data correctly, to learn from such data, and to use those learnings to achieve
 187 specific goals and tasks through flexible adaptation”. The definition encompasses
 188 three core aspects of AI systems. *Interpretation* of external data requires reason-
 189 ing, i.e., deriving conclusions from raw inputs using an internal representation
 190 of knowledge. *Learning* from data is the ability to change a system’s internal
 191 model based upon examples. *Adaptation* means the system can perform actions
 192 that change according to a change in the internal representation. In the follow-
 193 ing we describe the first two core aspects and their supporting technologies. The
 194 third aspect targets autonomous agents, such as robots, and game-playing (e.g.,

195 Silver et al., 2016), which are not relevant to the task of data integration and there-
196 fore are not addressed further.

197 **2.2.1. Knowledge representation and reasoning systems**

198 Allowing computer software to reason requires a way to represent and store
199 large amounts of knowledge, and systems able to query knowledge and rea-
200 son over it. One of the most mature approaches, backed by substantial com-
201 mercial and academic effort, is that of the Semantic Web as envisioned by
202 Berners-Lee and Hendler (2001). Under this conceptual model, knowledge graphs
203 (KG) have become a standard for representing facts. As their name suggests, KG
204 are a network-based representation, where entities and literals are nodes, and pred-
205 icates or relations are the edges.

206 **Example 3** *In Figure 3, a knowledge graph fragment presents our knowl-*
207 *edge about a data point from a dataset (Semina and Mikaelyan, 1994) stored*
208 *on PANGAEA. The dataset entity (Hydrolog...) is connected via the predicate*
209 *gl:hasProject to a literal describing it. The data point entity (Temp) is connected*
210 *via a predicate gl:hasDataset to the dataset entity describing the fact that the for-*
211 *mer is a component of the latter.*

212 In general-purpose knowledge graphs such as Wikipedia-based DBpedia
213 (Auer et al., 2007), entities may represent people, places, and abstract things, such
214 as events, while literals represent single pieces of information such as names,
215 titles, and dates. Ontologies provide a conceptualization of the domain (or do-
216 mains) described by the knowledge graph, adding entailment mechanisms such
217 as the ability to group entities into a class, create *same-as* links between en-
218 tities, equivalence relationships between classes, and denote predicates as sub-
219 properties. For example, both entities in the example above are connected via
220 *rdf:type* predicates to their ontological classes. These two entities and the predi-
221 cates prefixed with *gl:* are defined in the GeoLink base ontology (Krisnadhi et al.,
222 2015). The definition of an *rdf:type* is specified in the resource description
223 framework (RDF) and can be found at [https://www.w3.org/TR/rdf-](https://www.w3.org/TR/rdf-schema/)
224 [schema/](https://www.w3.org/TR/rdf-schema/). Querying information represented as a KG is often done using
225 SPARQL (Prud’hommeaux and Seaborne, 2008), a data retrieval language en-
226 hanced with semantic inference constructs.

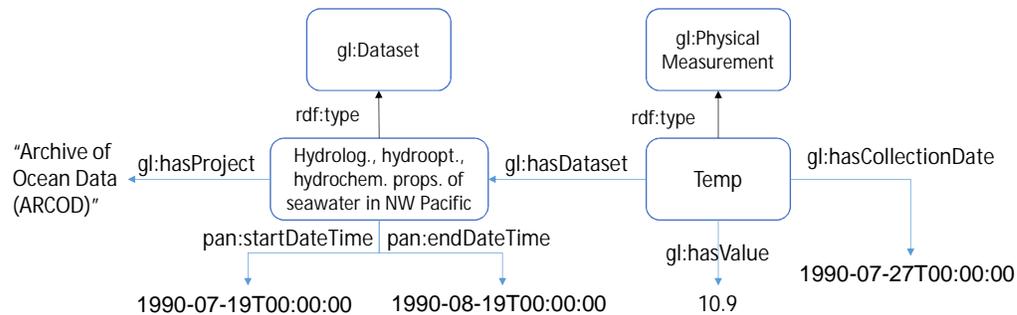


Figure 3. An example knowledge graph.

In the figure, a graph fragment with some of the data from Semina and Mikaelyan (1994) is presented in machine-readable manner by using well-defined ontological and schematic properties that have well-defined relations to other properties. These definitions and properties allow integrating these data with data from other datasets. Boxes represent entities, quoted strings are literals, and edges represent predicates that connect a subject (entity) to an object (entity/literal). Prefixes denote the ontology/schema in which the property/class are defined, with `rdf` denoting the resource description framework (RDF) schema (<https://www.w3.org/TR/rdf-schema/>), `gl` denoting the geolink ontology (<http://schema.geolink.org/>), and `pan` denoting the PANGAEA schema. The entity `Temp` represents a data point and is connected to its parent dataset via a `gl:hasDataset` predicate. The data point is connected to the collection time via a `gl:hasCollectionDate` predicate, and the dataset is connected to its temporal coverage through the predicates `pan:startDateTime` and `pan:endDateTime`. Both entities (i.e., data point and dataset) are connected to their ontological classes via an `rdf:type` predicate. The dataset entity is connected to a literal describing its project.

227 **2.2.2. Machine learning**

228 Endowing software with the ability to learn from examples has been studied exten-
229 sively over the past 60 years. ML has been used to automate tasks over the entire
230 expanse of the human endeavor from predicting relations in knowledge graphs
231 (see review by Nickel et al., 2016) to forecasting solar radiation (Voyant et al.,
232 2017). Machine learning techniques can be broadly divided into two types, *super-*
233 *vised* and *unsupervised* by the type of input provided to the learning algorithm.

234 Unsupervised learning techniques provide the learning algorithm with a large
235 collection of items sampled from the target population and some target metrics
236 to assess the quality of the task result, leaving the algorithm to attempt and op-
237 timize these quality criteria. Classic examples include clustering techniques such
238 as K-Means (Hartigan and Wong, 1979). The effectiveness and applicability of us-
239 ing unsupervised techniques to learn a representation have increased significantly
240 with the appearance of large amounts of user-generated content on the Internet.
241 For more details, see the seminal paper on the unreasonable effectiveness of data
242 by Halevy et al., (2009). A similar opportunity exists in oceanic sciences with the
243 increasing availability of large amounts of autonomously collected and remotely
244 sensed data (see Durden et al., 2017, for a review).

245 Supervised learning techniques require a (hopefully large) set of tagged ex-
246 amples. For example, to identify the semantic information conveyed by a set of
247 numbers representing the pixels in a picture, a supervised ML algorithm requires
248 a set of pictures labeled as *cats*, another labeled as *dogs*, etc. Similarly, to identify
249 people and places mentioned in a text, an ML model requires sentences where
250 they are clearly labeled as such. Given a metric to which the ML’s prediction can
251 be compared to the real tag, the ML algorithm can alter its internal representa-
252 tion to achieve better results on the task at hand. For example, using a quadratic
253 loss metric, calculated over the distance between the final result vector and the
254 expected one, is common in computer vision tasks. However, obtaining tagged
255 examples is often difficult and expensive, as it requires humans, often experts, to
256 tag the examples. Furthermore, one needs to obtain a set of examples which is
257 representative of the target task. More often than not, the examples on which ML-
258 models are trained are those for which gathering information is more convenient
259 than representative.

260 **2.2.3. Information extraction**

261 The ability of AI systems to obtain information from raw data relies upon three
262 fields of research. *Computer vision* (e.g., Krizhevsky et al., 2017) aims to ex-
263 tract meaning from images and video, (*textual*) *information extraction* focuses on

264 text (e.g., Martinez-Rodriguez et al., 2020), and *audio (speech) recognition* (e.g.,
 265 Hinton et al., 2012) converts sound into more meaningful information such as text
 266 and emotion markers (Schmidt and Kim, 2011).

267 2.3. AI in data integration

268 2.3.1. Ontology-based data integration and access

269 Taking advantage of the AI knowledge representation and inference mechanisms,
 270 *ontology-based data integration* (OBDI) uses ontologies to consolidate several
 271 heterogeneous sources into one source (see review by Ekaputra et al., 2017). For
 272 example, if the schema in one dataset contains the specific instrument (e.g.,
 273 CTD/Rosette) and in another the instrument type (e.g., Cast), we could use the
 274 *hasType* ontological construct to integrate them.

275 In many cases existing data sources are not linked to an ontology, rendering
 276 OBDI impossible. *Ontology-based data access* (OBDA) is an alternative model
 277 that provides access to the data layer through a declarative mapping between au-
 278 tonomous data layers and a domain-specified ontology (Xiao et al., 2018). A typ-
 279 ical development process of an OBDA system for a project that has a standard,
 280 non-ontological database will contain the following steps. (a) Create an ontol-
 281 ogy of domain-specific user knowledge. (b) Write mapping that connects (usually
 282 through SQL queries) the ontology to the project’s database. (c) Write a query us-
 283 ing ontology’s vocabulary as a semantic query language query, such as SPARQL.
 284 (d) Build an OBDA system framework that automatically rewrites the SPARQL
 285 query to the query language in which the project’s database operates.

286 2.3.2. Word embeddings

287 Early work in DI heavily relied upon measures such as Jaccard similarity (e.g.,
 288 He and Chang, 2006) and n-gram techniques (e.g., Do and Rahm, 2002) to ascer-
 289 tain if two strings are similar. However, syntactic methods ignore the semantics,
 290 or meaning, of words. Such techniques can find *plane* and *airplane* to be similar,
 291 but not *plane* and *aircraft*. To overcome this weakness, thesauruses such as Word-
 292 Net, and later Wikipedia, were introduced. However, these techniques required
 293 accurate spelling and were often baffled by technical terms and abbreviations.

294 The appearance of *word embeddings* has revolutionized the approach towards
 295 word, phrase, and sentence similarity. Word embedding was originally designed to
 296 convert text to the numerical required by DL techniques. The technique represents
 297 each word in the vocabulary with a d-dimensional vector of real numbers $w \in \mathbb{R}^d$.
 298 Word embedding has been extensively used in AI applications as an underlying

299 input representation that serves as a word dictionary and enables better capture of
300 the semantic meaning of the word (Levy et al., 2015). The following hypotheses
301 have been noted (Bolukbasi et al., 2016). (a) Vectors of words of similar meaning
302 tend to be closer. (b) The vector differences between vectors representing word
303 embeddings have been shown to represent relationships between words. A famous
304 example is the male/female relationship captured by the word2vec implementation
305 of word embedding, where Mikolov et al. (2013) showed that $\vec{King} - \vec{Man} +$
306 $\vec{Woman} \approx \vec{Queen}$.

307 Thus, a word would be embedded in a high-dimensional space as a vector,
308 and a sentence became a collection of such vectors. Word similarity now be-
309 comes a problem of vector similarity. Useful embeddings are those that place
310 similar words close to each other in this high-dimensional space. Embeddings
311 are learned from large collections of text, in an unsupervised manner. Thus, they
312 can be fine-tuned to a specific domain by retraining some of the embeddings
313 on a collection of domain-representative documents. Popularized by Word2Vec
314 (Mikolov et al., 2013), recent methods include GloVe (Pennington et al., 2014),
315 Flair (Akbiik et al., 2018), and BERT (Devlin et al., 2019). The latter two use
316 character-based embedding, which can also overcome spelling and abbreviation
317 issues.

318 2.3.3. Machine learning for data integration

319 The use of machine learning for schema matching had been pioneered by
320 Doan et al. (2000), followed by work by Gal and Sagi (2010). In both cases, ma-
321 chine learning was used to learn an ensemble model or method to combine the
322 results of multiple matchers by training the ensemble method on the results of
323 previous matching attempts. Sagi and Gal (2013) took this method one step fur-
324 ther by learning to adapt the ensemble weights according to the results of the ac-
325 tual matching performed at run-time. Thus, the features upon which their model
326 was trained were not the choice of matchers, but rather the structure and various
327 counting statistics of the match result. Recently, word embeddings were used to
328 enhance the effectiveness of schema matchers by Fernandez et al. (2018).

329 ML techniques have been used for entity resolution as well. Kenig and Gal
330 (2013) used an unsupervised ML technique called *maximal frequent item-sets*
331 (MFI) to learn the optimal clusters in which to search for duplicates. Sagi et al.
332 (2017) expanded upon this work by training an alternating decision tree model
333 (Freund and Mason, 1999) to classify pairs within the blocks to matched and un-
334 matched entities. Recent work, such as by Ebraheem et al. (2018), utilizes word

335 embedding to create semantically similar clusters as well as recommend matched
336 pairs. Data tamer (Gubanov et al., 2014) uses ML for entity consolidation by pre-
337 dicting which data item is most likely to be relevant.

338 **3. Data integration in ocean science**

339 In this section, we formalize the data integration process for oceanic datasets. Un-
340 der this formalization, we can compare similar tasks and examine tools employed
341 in support (or in relief) of the extensive manual labor otherwise required. After de-
342 scribing each step, we review current work in ocean science and list the remaining
343 gaps accompanied by specific directions for future work.

344 A data integration project can be described as having three major phases (Fig-
345 ure 4, top layer). In the *Discovery* phase, the list of possible candidate datasets for
346 the project is compiled. In the *Merge* phase, candidate datasets are harmonized
347 semantically, computationally, and geographically to form one large and coherent
348 dataset. In the *Evaluate/Correct* phase, the results are analyzed to assess quality,
349 coverage, and bias, and appropriate corrections are made to support assertions
350 made over the data.

351 In the following sections, we describe these phases in detail, further dividing
352 them into distinct steps. Although the integration process described holds whether
353 done manually or automated, we point out how the DI technologies described
354 in Section 2 can be used to automate the different steps, allowing to scale such
355 projects and integrate large amounts of data. Where appropriate, we describe how
356 AI technologies can in-turn support the DI processes. The bottom two layers of
357 Figure 4 summarize these supporting relationships.

358 *3.1. Discover*

359 Data discovery is the phase where candidate datasets are collected to fit a set of
360 study parameters. For example, Luo et al. (2012) searched for datasets contain-
361 ing sampling of marine N_2 (dinitrogen) fixing organisms. Similarly, Wang et al.
362 (2017) focused their efforts on geochemical data. The process of data discovery
363 can be divided into three distinct steps, namely, *search*, *link*, and *identify*, de-
364 scribed below.

365 **3.1.1. Search**

366 In the search step, a list of candidate research is collected. Search is performed on
367 repositories or through portals that provide access to multiple repositories, here-
368 after referred to as *data sources*. Data sources may contain either textual descrip-

369 tions of studies (i.e., scientific papers) or the datasets themselves. Google Scholar
370 is an example of a scientific portal to study descriptions, while PANGAEA is a
371 repository of datasets.

372 When searching for relevant research, users use search tools provided by the
373 data sources. These tools can be classified into one of three types of interfaces.
374 *Key word* queries comprise a sequence of terms of which at least one should be
375 present in the dataset for it to be returned in the results. *Ontological* queries rely
376 on well-defined ontological terms such as organism species or molecular com-
377 pounds, which the user specifies together with logical constraints and entailment
378 allowances to form a logical statement. Each candidate result must satisfy the log-
379 ical statement to be returned. *Parameter* queries rely on metadata associated with
380 the research, such as the publication date or the geographical location of the sam-
381 ples collected. Queries are formed by defining restrictions and combining them
382 using simple logical operators (and/or/not). To exemplify the difference between
383 ontological search and parameter search, consider the following.

384 **Example 4** *A researcher is interested in datasets containing measurements of*
385 *phytoplankton biomass, among other parameters. In a parameter search, that re-*
386 *searcher would be required to search for all possible subgroups and types of phy-*
387 *toplankton, such as diatoms, Fragilariophyceae, and Coscinodiscophyceae, and*
388 *then collate the results. In an ontological search, the researcher can simply ask*
389 *for all diatoms and specify that they wish for all sub-species as well, then receive*
390 *all datasets containing the biomass of a species present in the taxonomic tree un-*
391 *der diatoms. However, to support such a search, each parameter defined over a*
392 *dataset needs to be aligned correctly with a comprehensive ontology, a task that*
393 *is daunting when done retrospectively over large collections of datasets.*

394 Table 2 provides a partial list of data sources, oceanic research portals and
395 repositories current to January 2020, their type (R stands for Repository and P
396 for Portal), and the extent to which they support the search tools described above
397 (all data Sources listed provide key-word search). A notable omission from this
398 list is the set of commercial cloud services participating in NOAA's Big Data
399 Project (National Oceanic and Atmospheric Administration, 2020). Access to this
400 data source is rudimentary, and the number of datasets provided is limited.

401 Taxonomies are widely used in the ocean sciences (Claramunt et al., 2017).
402 Some examples are *World Register of Marine Species* (WoRMS Editorial Board,
403 2020) that holds a detailed taxonomy of marine species, *AlgaeBase*
404 (Guiry and Guiry, 2020), a global algal database, and *FishBase* (Froese and Pauly,

405 2020). An *ontology* is an explicit specification of a conceptualization that defines
406 the terms in the domain and relations among them (Gruber, 1995).

407 All ontologies use some form of *vocabularies* in order to express terms and
408 specify their meanings (Uschold, 1998). Similarly to taxonomies, they adopt a
409 classification structure. However, ontologies add properties for each class and a
410 set of axioms and rules that allow reasoning and full domain conceptualization
411 (Zeng, 2008). Leadbetter et al. (2010) provide a systematic review of ontologies
412 for the maritime domain. A few notable mentions include the NASA Semantic
413 Web for Earth and Environmental Terminology (SWEET; Ashish (2005)), which
414 hosted over 6000 concepts in 200 separate ontologies as recently as 2018, but
415 since 2019 has been removed from public access. MarineTLO is a top-level on-
416 tology for the maritime domain (Tzitzikas et al., 2013) that contains informa-
417 tion about marine species, ecosystems, and fishers. Significant among these ef-
418 forts is OceanLink/GeoLink, a large-scale project that aims to improve discov-
419 ery, access, and integration (Figure 4) of interdisciplinary data in the oceanog-
420 raphic domain (Narock et al., 2014). The ongoing project enables the discov-
421 ery of integrated data from multiple repositories by creating an integrated knowl-
422 edge discovery framework on top of those repositories. The project utilizes se-
423 mantic web technologies, particularly ontology design patterns (ODPs; Gangemi
424 (2005)) and a SPARQL endpoint (accessible at data.geolink.org/sparql) for se-
425 mantic querying. Additional repositories supporting OBDA through a SPARQL
426 endpoint are the European Directory of Marine Environmental Data (EDMED) (at
427 <https://edmed.seadatanet.org/sparql/>), and the British oceanog-
428 raphic data centre, NERC SPARQL endpoint (at [http://vocab.nerc.ac.
429 uk/sparql/](http://vocab.nerc.ac.uk/sparql/)).

430 Although GeoLink's ontologies provide extensive coverage of the domain,
431 they are far from complete. In some cases, publishing a repository's data in Ge-
432 oLink is not possible due to missing concepts or a required but tedious schema-
433 mapping process that the authors do not wish to undertake. In those cases, the
434 remainder of the data not described by the GeoLink ontologies is published ac-
435 cording to the provider's own schema (Krisnadhi et al., 2015). Specifically, some
436 of the more fine-grained patterns are not fully described. For instance, in the ma-
437 rine biology domain, integrating data according to taxonomy can be very useful.
438 Similarly, for measurements of plankton data such as biomass, integrating data
439 according to plankton group size or kind can be beneficial. Such a taxonomic
440 relation exists in the MarineTLO ontology and in WORMS but is missing in Ge-
441 oLink. Another example is the lack of ontological representation of ocean basins
442 and seas such as in SeaVoX (Claus et al., 2014). The GeoLink class *Place* can be

443 related to a *PlaceType*= 'ocean' but no deeper hierarchical representation is sup-
444 ported. For example, if the discussed place is set to 'The Red Sea' and some other
445 data point is given with the place set to 'Gulf of Eilat' (a part of the Red Sea)
446 the correct integration could not be made with Geolink. Even if the ontological
447 issues are resolved, realigning existing data with Geolink, or a combination of the
448 existing ontologies, would require an extensive mapping effort that would benefit
449 from AI-supported schema matching technologies.

450 Thus, scaling the search process by using OBDA would allow the collection
451 of a large number of datasets already aligned by the domain ontology over the
452 parameters used to perform the search step. However, using OBDA requires the
453 domain ontology to cover all aspects of the data to be integrated, and all datasets
454 in the repository/portal to be completely aligned with the ontology. As detailed
455 above, current repositories and data portals mostly use taxonomies rather than
456 ontologies, combining parameter and keyword search. Existing domain ontolo-
457 gies have limited coverage and cross-alignment. In the absence of perfect OBDA
458 systems, the *merge* phase is required to integrate the different datasets with their
459 mismatched schemas and data descriptions.

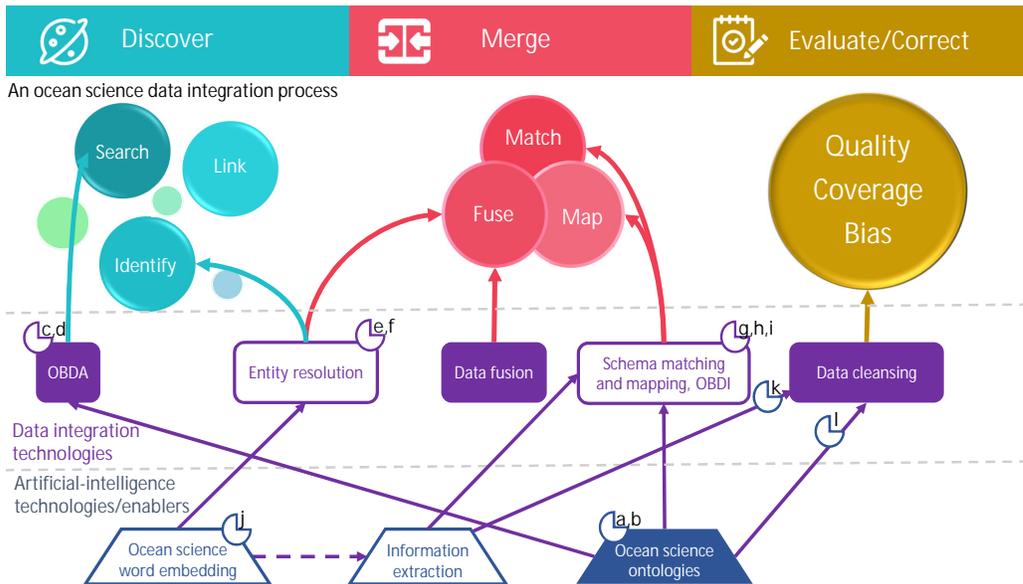


Figure 4. The three phases of the data integration process, and their application in ocean science.

The top layer describes the process: in the *discover* phase, a list of candidate datasets with possible relevancy to the project is compiled; in the *merge* phase, candidate datasets are harmonized semantically, computationally, and geographically to form one large and coherent dataset; in the *evaluate/correct* phase, an analysis of the resulting dataset is performed to assess quality, coverage and bias, followed by appropriate corrections that are made to support assertions made over the data. The middle layer shows how data integration technologies support the process. OBDA and OBD stand for *ontology-based data access* (A) and *integration* (I) respectively. The bottom layer contains three AI technologies/enablers that support the data integration technologies. Full-colored rectangles and trapezoids represent technologies/enablers in current use. Outline-colored-only shapes represent technologies and enablers that are not currently in use in ocean science data integration. Additional gaps are listed as lower-case letters corresponding to the gaps listed in Table 3.

Table 2. Examples of oceanic data sources.

Data source	Type^a	Content type	Ontological support	Searchable parameters (excl. key words)
ARGO	R	Float	No	Date, geo-coordinates
BCO-DMO	R	Underway, cast, float	No	Date, geo-coordinates
COPERNICUS	P	2D/3D images, cast, float	No	Date, geo-region, parameter name
EDMED	R	Underway, cast, float	Yes	Date, geo-region, geo-coordinates, parameter (ontology), instrument (ontology)
Global DMS	R	Underway	No	Date, geo-coordinates
Google dataset search	P	All	No	None
IsraMar	R	Cast	No	Date, geo-coordinates, parameter name
NCEI LAS	R	Cast, underway, 2D image, radar, float	No	Date, geo-coordinates
PANGAEA	R	Cast, underway, float	No	Date, geo-coordinates, geo-region, instrument
SeaBass	R	Cast, 2D image	No	Date, geo-coordinates, instrument
World ocean database	R	Cast, underway, 2D image, radar, float	Yes	Date, geo-coordinates, instrument, parameter name, bio-species (ontology)
Data One	P	All	Yes	Date, Geo-coordinates, instrument, parameter name, bio-species (taxonomy)

^a R: data repositories. P: portals. Portals provide access to data from multiple repositories.

460 **3.1.2. Link**

461 The linking process entails connecting between studies and their datasets (and
462 vice versa) and between datasets, which are derived from one or more other
463 datasets. The prevalence of object identifiers such as DOI, coupled with the in-
464 creasing tendency of authors and publishers to provide publicly available datasets
465 together with submitted papers, has made this process easier. However, the link-
466 ing process is still a largely manual process where researchers piece together
467 the papers describing the data and vice versa. Furthermore, the linking process
468 may require a finer resolution, as the following story published by Data One
469 (Data Observation Network for Earth, 2020) exemplifies.

470 “A third dataset looked particularly promising for use in a global
471 study, but its PI had neglected to include units of measurement in
472 the dataset. Unwilling to give up on a potentially great contribution,
473 Eileen decided to do some detective work and pull up the associated
474 publication, looking for any clues that might lead to a breakthrough.
475 At long last, Eileen found a single table referencing the units for a par-
476 ticular column of data. With the units finally established, she worked
477 backwards to make sense of the data – but at a cost of several hours’
478 work.”

479 Thus, even though the researcher had succeeded in linking the dataset to its
480 corresponding publication, more refined work was needed to link specific param-
481 eters to their descriptions. This refined linkage can be delayed until the merge phase
482 where the extended data descriptions can be used to better align the schemas of
483 the integrated datasets’ with the domain ontology.

484 **3.1.3. Identify**

485 Even with the existence of DOI, in many cases, the same data may appear in
486 several datasets by being used for several studies. Thus, researchers are required to
487 meticulously read the data collection procedures of every study used to make sure
488 that their data do not contain duplicate measurements and identify each dataset
489 or even data point in a unique manner. The implicit danger of duplicates is that
490 they can create an inherent bias in the results towards duplicated data. In oceanic
491 repository integration, this process is further complicated by the fact that some
492 records represent a collection of datasets that previously may have been published
493 separately as well.

494 Thus, DOIs provide grounding of datasets to fixed, reliable repository men-
495 tions, and can be used for citation and referencing purposes. However, they do

496 little to resolve issues such as data overlap, republication, and bundling that
497 may manifest themselves when combining several datasets. Resolving duplicate
498 datasets and overlapping data points using entity resolution (see Section 2.1) is
499 an obvious use of AI-supported DI tools. As entity resolution tools rely on simi-
500 larity comparisons, they would also be benefited by ocean-science-specific word
501 embedding to allow semantic comparison.

502 3.2. *Merge*

503 Once a collection of datasets has been assembled, the *merge* phase can com-
504 mence. To facilitate this process, one must create a mediated schema to which
505 all other datasets are matched and subsequently mapped or use an ontology to
506 which the datasets' schemas are mapped to facilitate OBDI. We divide this phase
507 into three distinct steps, described in detail below. In the *match* step, correspon-
508 dences are found between each attribute in every dataset's schema and the medi-
509 ated schema/ontology. In the *map* phase, a function mapping from the semantics
510 of the source dataset's schema to the mediated schema is constructed. In the *fuse*
511 step, some datasets are interpolated over space/time to create a continuous and
512 uniform space of measurements.

513 3.2.1. **Match**

514 In the match step, researchers align the different attributes/parameters in the
515 dataset's schema with the mediated schema/ontology. To do so, the researcher
516 must often consult the data descriptions of each parameter, which are either listed
517 with the dataset in the source repository or described as part of the methods sec-
518 tion of the accompanying paper. If an exact match cannot be found, the researcher
519 must decide whether to disqualify the parameter or even the whole dataset from
520 inclusion in the study or extend the mediated schema/ontology to accommodate
521 the new dataset.

522 A wealth of literature and tools exist in the general database and knowledge-
523 base domains to facilitate schema matching and ontology alignment. Among these
524 are the use of acronym expansion (e.g., Sorrentino et al., 2010), a corpus of pre-
525 viously discovered correspondences (e.g., Madhavan et al., 2005), and instance
526 information (e.g., Chen et al., 2018). However, to the best of our knowledge, none
527 of these were applied to match ocean science dataset schemas, neither pair-wise
528 nor to mediated schemas or ontologies. Zhou et al. (2018) proposed a complex
529 real-world ontology alignment benchmark made on two separate GeoLink dataset
530 ontologies. However, even this unique example attempts to automate ontology

531 alignment and not automatically match dataset schemas against these ontolo-
532 gies. Furthermore, none of the existing automated schema matching and mapping
533 tools is interoperable with the common ocean science meta-data formats. Schema
534 matching can be supported further by AI-based information extraction technolo-
535 gies, such as described in Section 2.2.3, by extracting data descriptions from the
536 research papers linked to the datasets. These data descriptions can be used to im-
537 prove schema matching performance, thus utilizing this unique aspect of ocean
538 science datasets.

539 **3.2.2. Map**

540 In some cases, the semantics of the data in one source are slightly different from
541 that of the mediated schema/ontology. For example, a dataset may contain two
542 fields, one representing the latitude and another the longitude, while in the me-
543 diated schema, there exists a single *coordinates* field that combines the two. In
544 other cases, the mediated schema may contain a field that represents a calculation
545 performed over raw data, or the units of measurement may differ between sources.
546 All of these examples, and other semantic differences, require a mapping phase
547 where conversion functions are generated to facilitate data integration according
548 to correspondences found in the matching step. Even more mundane, but crucial
549 is the need to map from the source format to that of the central repository used to
550 collect the data from the different datasets. For example, the data may be received
551 in XML format and the repository stored in a relational database, requiring for-
552 mat conversion between the two. The use of OBDI facilitates conversion between
553 fields of different datasets by using the encoded conversion logic within the ontol-
554 ogy. Thus, for example, the concept of *Celsius*, can be linked to the concept of
555 *Fahrenheit* by a relation containing a specific bi-directional conversion function.

556 Together with the match step, there is a substantial need for golden-standard
557 tasks and structured benchmarks for ocean science schema-matching and map-
558 ping tasks to enable the development and training of automated matching tools
559 utilizing the existing ontologies and vocabularies. Word-embedding-based tools
560 are highly dependent on the domain from which the text used to generate the em-
561 bedding was collected. Currently absent, a word embedding for the ocean science
562 domain would be an important enabler for AI-based DI tools (see Section 4). The
563 same embedding could be used to enhance information extraction tools to sup-
564 plement schema matching and mapping processes over datasets with information
565 from their linked papers. As a foundational enabler, providing schema interoper-
566 ability between the common ocean science data formats and those used by schema
567 matching and mapping tools would open up a plethora of options for practitioners

568 to use.

569 **3.2.3. Fuse**

570 In this step, researchers need to mitigate problems that emanate from differences
571 in spatio-temporal resolution between the datasets. Thus, one dataset may include
572 measurements of a 50-m depth in increments of 1 m, while another in increments
573 of 10 cm. Decisions must be made on whether to aggregate upwards to lower
574 resolutions, omit incompatible resolutions or interpolate the data to align the res-
575 olutions, or fill out missing data in some areas (e.g., in Kaplan and Lekien, 2007,
576 due to faulty sensors). As previously mentioned, we leave the review and critical
577 analysis of existing work in data fusion to future work.

578 In addition to spatio-temporal fusion, this step entails an additional effort of
579 resolving duplicate and overlapping data points. While overlapping and duplicate
580 datasets could possibly be identified at the *identify* step, identifying these cases at
581 the datapoint level requires all fields to be aligned by the match and map steps.
582 Here, again, we can use entity resolution to automate this task (see Example 2).

583 *3.3. Evaluate and correct*

584 After, or sometimes during, the data integration process, researchers must evalu-
585 ate the integrated dataset to facilitate inclusion/exclusion decisions and to report
586 quality and descriptive measures upon publication. The evaluation process often
587 addresses one or more of the following issues.

588 **3.3.1. Quality**

589 Detecting data errors is often done using non-specific numerical and statistical
590 tools; for example, by excluding all outliers, defined as values over two standard
591 deviations from the mean. This step can be mostly aligned with the existing DI
592 process of *data cleansing* (see Section 2.1). To identify, quantify, and possibly
593 correct errors in data via interpolation, techniques appropriate to the data type
594 (e.g., Gupta et al., 2014) should be used. Here, we refrain from performing a de-
595 tailed review of the extent of AI used in these processes over ocean science data
596 in the interest of brevity and focus.

597 A non-generic approach that could provide more accurate results can be ob-
598 tained by reasoning over accumulated knowledge tied to the domain ontologies.
599 For example, O'Brien et al. (2013) needed to remove individual samples of occol-
600 ithophore (a type of plankton) where the species was reported as *Thoracosphaera*
601 *heimii*, as this species was reclassified out of the coccolithophore family after

602 the original data were collected. This removal of misclassified samples could be
603 done automatically by defining a logical rule over the global ontology. Further-
604 more, among the tools that can support a researcher in the process of evaluating
605 the data quality of a given dataset, information extraction can provide substan-
606 tial assistance. For example, information extraction tools can be used to extract
607 and categorize quality control processes and pre-processing techniques used in
608 a specific dataset and a collection of datasets from the scientific text describing
609 them. Once extracted, this information can be attributed to the dataset, allowing
610 researchers to employ data cleansing methods and filter out less trustworthy pro-
611 cesses or, conversely, to select only those data points on which the required type
612 of pre-processing was performed.

613 **3.3.2. Coverage and bias**

614 An important tool in the evaluation of result validity and relevance is the analysis
615 of coverage and bias. Data are collected in different geographical regions, depths,
616 and seasons, and using different instruments. When presenting results, one must
617 either correct them for inherent biases, exclude under-represented partitions, or
618 provide a list of caveats and analyses regarding the coverage and bias with respect
619 to the general distribution over each dimension (geographical/temporal/other).
620 The ability of an ocean scientist to make use of an AI-based integrated dataset
621 strongly depends on accurate representation of possible biases and uncertainties
622 associated with the DI process. This point is emphasized for the case of climate
623 science studies, where uncertainties result from a wide range of sources, as a lim-
624 ited number of available measurements, especially for rare events (IPCC, 2014).

625 Existing portals/repositories provide mechanisms to filter by time/geo-
626 location or map a collection of datasets over a world map. These mechanisms
627 allow researchers to assess the coverage of their collection of datasets if they are
628 from the same portal/repository. Evaluating coverage and bias over other dimen-
629 sions, such as instruments used and bio-diversity, is dependent on the ability to
630 perform OBDA, the coverage of the OBDA's ontology, and the extent of informa-
631 tion extracted from the scientific description and aligned with the ontology.

632 *3.4. Summary*

633 Figure 4 presents an overview of how DI technologies (in purple/purple outline,
634 middle layer) could support and scale the different steps and phases of the ocean
635 science data integration process. However, to make these technologies work, some
636 AI technologies and enablers are needed. These are listed in the bottom layer

637 of the figure as trapezoids and are connected to the DI technologies which they
638 support. Ontology-based technology features heavily, as it effectively combines
639 the wealth of accumulated knowledge of the oceanic domain with AI-supported
640 DI technologies. DI technologies and AI technologies/enablers that are missing
641 today are drawn with a white background.

642 Table 3 presents a list of existing and missing enablers for DI in ocean sci-
643 ence. Some of these enablers are presented in the figure, while others enable the
644 processes in the figure. The gaps in the table are annotated with lower-case letters
645 that are repeated in Figure 4 where they are positioned on the DI technology they
646 enable, on the AI technology they enable, or on the support a specific AI tech-
647 nology provides to a DI technique. Note that while the technologies and enablers
648 reviewed in Table 3 are listed by phase, some of them support multiple phases. For
649 example, *entity resolution* is a DI technology that can be used to identify duplicate
650 datasets prior to their integration in the *identify* step and to identify duplicate data
651 points in a merged dataset as part of the *fuse* step.

Table 3. Missing and existing AI enablers for DI in ocean science

Phase	Existing enablers	Remaining gaps
Discover	(1) Several ocean science ontologies. (2) OBDA to major dataset repositories. (3) Extensive use of DOI.	(a) Incomplete conceptual coverage of existing ontologies. (b) Incomplete conceptual alignment between ontologies. (c) Alignment of historical datasets with existing ontologies. (d) AI-based tools for creators to align their schemas with existing ontologies.
Merge	(4) An ocean science ontology alignment benchmark.	(e) Entity resolution oceanographic benchmarks for both dataset and data point levels. (f) Entity resolution tools utilizing ocean science word embeddings. (g) Ocean data format interoperability with existing tools. (h) Schema matching/mapping oceanographic benchmarks. (i) Matching and mapping tool utilizing semantics encoded in existing vocabularies and ontologies. (j) Word embedding for ocean science domains.
Evaluate	(5) Existing work on data cleansing/anomaly detection. Not reviewed in detail. (6) Geolocation mapping in data-portals	(k) Annotated datasets, tools, and benchmarks for extracting data quality and preprocessing descriptions from scientific text (l) Extension and refinement of oceanographic ontologies with respect to coverage, bias and quality queries.

652 **4. Empirical evaluation: the impact of AI infrastructure**

653 In the following section, we provide some empirical evidence to the necessity of
654 creating the AI infrastructure required to support DI efforts in ocean science. As
655 described in the previous sections, both AI-supported entity resolution tasks in
656 the *discovery* phase and schema matching tasks in the *merge* phase could bene-
657 fit from adding relevant information from unstructured sources accompanying the
658 data. In Example 1, the fact that the Nitrate field represented the sum of nitrate and
659 nitrite was mentioned in the column comments. The ability to retrieve this infor-
660 mation from the comment, codify and align it with a domain ontology, relies on
661 AI-software being able to recognize domain-specific information in unstructured
662 text. Domain-specific datasets, benchmarks, and word embeddings are needed to
663 bridge this gap (see Table 3). To exemplify the potential benefits of having this
664 infrastructure in place, we train a state-of-the-art information extraction system
665 on ocean science data descriptions and report on the performance gains on an
666 information extraction task.

667 *4.1. The task: extracting data descriptions using information* 668 *extraction techniques*

669 A standard information extraction task, named entity extraction (NER) aims to
670 find entity mentions in unstructured text and map them into predefined classes.
671 These entities can then be used to enrich automated data integration tasks such as
672 schema matching and mapping. The classes a NER is seeking in the text can vary
673 based on the requirement of a specific assignment. The most widely used classes
674 are person, location, organization, and date (Jiang et al., 2016). For instance, a
675 NER system trained to detect person, location, and organization when receiving
676 the following text as input: "John Doe lives in New York City and works in the
677 New York stock exchange," should identify the following named entities as out-
678 put, where the named entity is denoted between brackets and the class between
679 parentheses. [John Doe](person) lives in [New York City](location) and works
680 in the [New York stock exchange](organization). An ocean science DI application
681 would need to identify entities such as a measured variable (temperature, salinity),
682 units (degrees, dbar), and devices (CTD, sonar, plankton counters).

683 *4.2. Datasets*

684 **4.2.1. An oceanic science entity extraction dataset**

685 To the best of our knowledge, no gold-standard annotated documents are
686 freely available for the oceanic domain. Therefore, we created a small dataset

687 to provide initial support to our claim for the need for an extensive stan-
 688 dard to train and evaluate tools against. We retrieved 30 documents contain-
 689 ing data descriptions from three data repositories: PANGEA (2020), BCO-
 690 DMO (Biological and Chemical Oceanography Data Management Office, 2020),
 691 and the European directory of marine environmental data (EDMED, 2020). Each
 692 token (usually a single word) was annotated in the IOB2 format using a standard
 693 NER annotation tool named TALEN (Mayhew and Roth, 2018). The IOB2 format
 694 is a tagging format designed for the NER task. The *B*- prefix before a class name
 695 is used to indicate that the token is at the beginning of a chunk, the *I*- prefix before
 696 a class indicates that the token is inside a chunk, while *O* represents a token that
 697 is not inside of any chunk. Figure 5 shows an example of the IOB2 format used
 698 to annotate a data description document retrieved from EDMED. Our test data
 699 contain 1,256 sentences and 7,848 total tokens with an average of 262 tokens per
 700 document. We found 2,193 entities divided into 11 classes averaging 75.6 enti-
 701 ties per document with an average length of 2.17 tokens per entity. The dataset is
 702 available online (Bar, 2020a).

```

Environmental O
modification O
caused O
by O
aquaculture O
along O
the O
Portuguese B-GeoRegion
continental I-GeoRegion
coast I-GeoRegion

```

Figure 5. An example of the IOB2 annotation.

In this figure the IOB2 annotation is used to identify a GeoRegion within a data description document retrieved from EDMED. Tokens marked with an O are not part of any entity. The token marked with B-GeoRegion begins the entity. The rest of the entity’s tokens are marked with I-GeoRegion.

703 4.2.2. An oceanographic text dataset

704 Word embeddings are created using a large text corpus. To test the hypothesis that
 705 specific word embedding could improve NER algorithms on the task of identifying

706 oceanic entities in texts, we trained custom word vectors. Our training method is
707 constructed based on the following steps. (a) Collect a large set of oceanographic
708 papers. (b) Extract raw text from the collected oceanographic papers. (c) Train
709 word embeddings based on the text corpus.

710 Due to overlapping terms from the oceanic domain in other closely related sci-
711 entific domains such as earth science or biomedical science, we collected papers
712 that were published in known oceanographic journals. We used the Crossref API
713 (Lammey, 2015) to search for the DOIs of papers that appeared in oceanographic
714 journals, such as *Ocean Science*, *Frontiers in Marine Science*, and *Aquatic Biol-*
715 *ogy*.

716 After acquiring the relevant DOIs, we implemented a web crawler that
717 searched for the full-text PDF version of the papers in several public reposi-
718 tories. The crawler mined 30,000 oceanic papers. We used the *Science Parse*
719 (Clark and Divvala, 2015) open-source Java library to extract data from the pa-
720 pers. We extracted the title, abstract, and content section parts of the documents
721 (references were excluded) into a JSON format. The raw text from the JSON file
722 contained over 175 million tokens. This dataset is available online as well (Bar,
723 2020b).

724 4.3. *Methods*

725 The NER algorithm is a supervised ML model that is trained on annotated docu-
726 ments to recognize patterns identifying a token or set of tokens as a named entity
727 and to which class it most likely belongs. For example, after seeing a large num-
728 ber of documents where the tokens next to the word *lives* describe a person (e.g.,
729 John Doe lives in), the ML model learns to classify these tokens as people. Us-
730 ing word embeddings to represent the documents on which the algorithm trains
731 allows it to generalize its learned model so that similar words such as *resides* and
732 *works* would be recognized as well. Furthermore, the token *John* itself is embed-
733 ded into the vector space such that other people’s names will be situated close to
734 it. As described in Section 2, generating word embeddings is an unsupervised ML
735 technique based on the co-occurrence of words in a very large text corpus.

736 In this evaluation, we use the Flair NER algorithm (Akbik et al., 2018), which
737 is based on a word embedding technique as well. Unlike other models, the
738 model employs character level tokenization rather than word-level tokenization.
739 A sentence is converted to a sequence of characters, and through a language
740 model, the algorithm learns the word representation. Flair uses a stacked em-
741 bedding approach. The algorithm’s character language model vector is concate-

742 nated with GloVe’s word embeddings (Pennington et al., 2014) to form the final
743 word vectors, thus leading to a better result. Flair produced state-of-the-art F1-
744 scores on the CoNLL-03 general-purpose dataset collected from newspaper arti-
745 cles (Sang and De Meulder, 2003).

746 To adapt Flair and its NER algorithm to the oceanic domain, one can both
747 retrain it (using supervised ML) on the classes of this domain and fine-tune the
748 underlying word embeddings (using unsupervised ML) to reflect semantic rela-
749 tions in this domain better. In the following, we demonstrate both improvements.

750 **4.3.1. Improving the Flair NER by retraining on an ocean science tagged** 751 **dataset**

752 Training was performed on a Gigabyte Technology server with an Intel i7-7700
753 8 core CPU, 64GB RAM, and Gigabyte GTX 1070 GPU running the Ubuntu
754 16.04.6 operating system. The empirical evaluation was performed using Flair
755 version 0.4.1 (Zalando Research, 2019) running on python version 3.6.8, deployed
756 as part of the Anaconda data science platform (Anaconda, 2020). We split our
757 annotated dataset randomly into a training set comprised of 80% of the documents
758 and a test set comprised of the remaining 20%. We then proceeded to train the Flair
759 algorithm on the training set and test both the original Flair NER model and our
760 retrained one on the test set.

761 **4.3.2. Creating ocean science word embeddings**

762 We utilized the oceanographic text corpus for training two new word embed-
763 dings. Word2Vec (Mikolov et al., 2013) with word-level embeddings and Flair’s
764 character-based forward and backwards embeddings (from now on, CFBF). The
765 word-level embeddings were implemented using the *emphGensim* Python library
766 (Řehůřek and Sojka, 2010) and the CFBF embeddings with Flair. One of the
767 known connections in oceanographic research is between a measured variable
768 and its measured units. Although often a variable can be measured using different
769 units, some notations are very common in the scientific literature. Similar to the
770 King-Queen relationship stated by Mikolov et al. (2013) on general-purpose text,
771 the oceanographic trained models were able to conclude the relationships in Fig-
772 ure 6. Recall that in general text, the vector representation of the word *king* was
773 found to relate to the vector representing *queen* in the same manner as the vector
774 *man* relates to the vector representing *woman*. After reviewing the ocean science
775 research papers, the unsupervised algorithm, with no input from a domain expert,
776 created an embedding model where, e.g., *m/s* relates to *speed* in the same manner
777 that *PSU* (practical salinity units) relates to *salinity*. Note, that the fact that PSU

$$\begin{aligned} \overrightarrow{\text{temperature}} - \overrightarrow{\text{salinity}} + \overrightarrow{\text{PSU}} &\approx \overrightarrow{\text{degrees}} \\ \overrightarrow{\text{speed}} - \overrightarrow{\text{salinity}} + \overrightarrow{\text{PSU}} &\approx \overrightarrow{\text{m/s}} \\ \overrightarrow{\text{pressure}} - \overrightarrow{\text{salinity}} + \overrightarrow{\text{PSU}} &\approx \overrightarrow{\text{MPa}} \end{aligned}$$

Figure 6. Variable-unit analogies.

The figure shows semantic relations between oceanic variables and their associated units, as found by the word embedding algorithm, with no intervention of a domain expert. Note that although salinity is a unitless variable, it was associated with PSU (practical salinity units) by the unsupervised algorithm. The relations can be read as follows. “*temperature* relates to *degrees* as *salinity* relates to *PSU*”

778 has since been retired is unknown to the embedding algorithm, as it was trained
779 on papers using this unit. Rather, this domain knowledge should be coded into an
780 ontology to ensure that data from papers using PSU, can be handled appropriately
781 when integrated with more modern datasets.

782 We trained the Flair algorithm with the same 80%-20% train-test split to detect
783 data descriptions from unstructured data, where the word embeddings served as
784 features for the NER algorithm. We ran the following stacked embeddings mod-
785 els: (a) GloVe and Flair embeddings trained on a general-purpose text that served
786 as a baseline; (b) Word2Vec oceanographic model; (c) Flair’s CFBF embeddings
787 trained on an oceanographic corpus; (d) Stacked embeddings model that was com-
788 piled of (b) (c) embeddings; and finally (e) Stacked embeddings model of (a) (b)
789 (c).

790 4.4. Evaluation measures

791 Several evaluation metrics have been offered to assess the efficacy of a NER sys-
792 tem, where the most commonly used are based on the exact-match evaluation.
793 A named entity that has been proposed by a NER system is considered correct
794 only if there is an exact match of both entity boundaries and class (i.e., all tokens
795 that should belong to the entity are correctly marked and assigned). However, the
796 ML model we use in the first evaluation was not designed to detect ocean science
797 classes (e.g., measure variable). As a result, we seek an exact boundary match with
798 no consideration of the entity type. For example, if the NER system can detect the
799 ‘Mediterranean Sea’ as a named entity, it will be considered a match regardless of
800 the class (*location* in this example). If for the same sentence, the system will only

801 detect 'sea' as a named entity, it will be considered a false match. In the second
 802 evaluation, we train all models to detect the specific class as well as extract the
 803 named entity and therefore seek an exact match of both boundary and entity type.

804 The measures precision, recall, and F1-score are arguably the most commonly
 805 used to aggregate and quantify the number of exact matches detected by a NER
 806 system. *Precision* is the fraction of true instances of the total number of instances
 807 predicted by the NER system as positive, while *recall* is the fraction of true in-
 808 stances predicted by the NER system of the total true instances in the dataset.
 809 *F1-score* is the harmonic mean of precision and recall. Their formal definitions
 810 are as follows.

Definition 1 (NER evaluation measures) *Let predicted positive (PP) be the set of named entities predicted as such by a NER algorithm. Let actual positive (AP) be the set of named entities that actually exist in the task. Let true positive (TP) be the intersection between these sets, i.e., those named entities that both exist in the task and were predicted by the NER algorithm, then Precision, Recall, and F1-score are defined as follows.*

$$Precision = \frac{TP}{PP} \quad (1)$$

$$Recall = \frac{TP}{AP} \quad (2)$$

$$F1 = 2 * \frac{Precision * Recall}{Precision + Recall} \quad (3)$$

811 4.5. Results

812 The result of the first evaluation can be seen in Table 4. The F1 score of the
 813 original flair model on oceanic data is only 0.068. Training the same flair model
 814 on an oceanic dataset results in an F1 score of 0.738. The results of the second
 815 evaluation can be seen in Table 5. The best model was the stacked embeddings
 816 model that reached an F1 score of 0.679 on unstructured metadata. We remind
 817 the reader that in the first task, we require only a boundary match, while in the
 818 second, we require both boundary and class to be correct, making it substantially
 819 more difficult.

Table 4. Performance of data description extraction using embeddings trained on general versus ocean science text.

Measurement ^a	Flair NER using news-trained embeddings	Flair NER using ocean-science-trained embeddings
Precision	0.221	0.746
Recall	0.040	0.731
F1 score	0.068	0.738

^a In this task, a true positive result entails identifying a named entity regardless of its class.

820 4.6. Discussion

821 The unmodified Flair model used in this evaluation scored a 0.932 F1 score on
 822 newswire text (Akbik et al., 2018). The same algorithm fails miserably on our
 823 task. The results of the retrained model can be considered as an immense im-
 824 provement but still far from state-of-the-art results achieved on NER tasks in
 825 other domains. This result is expected due to the small number of training ex-
 826 amples available to the supervised training algorithm. The result also highlights
 827 the need for an extensive, well defined, annotated dataset to train ML models
 828 over oceanic sciences tasks. Furthermore, the classes used to extract information
 829 should be carefully aligned with ocean science domain ontologies if they are to be
 830 used in conjunction with schema matching tools.

831 The oceanic embeddings allow Flair to boost its results on the harder bound-
 832 ary+class task from an F-1 of 0.415 to 0.679 for the best model. Here, too, a much
 833 more substantial increase is expected should we increase the amount of training
 834 data. Alternatively, we could use transfer learning from models trained on related
 835 datasets, such as scientific papers in general. Although 175 million tokens may
 836 sound impressive, the standard GloVe vectors used in general-purpose tasks are
 837 trained over 6 to 840 billion tokens (see Pennington et al., 2020, for examples).

838 5. Conclusions and future work

839 The study of the oceans relies on the extensive collection of physical, chemical,
 840 and biological data from various locations around the globe. Over the last cen-
 841 tury, numerous measurements have been performed continuously, resulting in the
 842 creation of an increasingly large amount of oceanic data. One of the significant

Table 5. Comparative performance of Flair NER using oceanic word embeddings as features.

Embeddings method	P^a	R	F1
Flair + GloVe (General-purpose)	0.547	0.335	0.415
Oceanic Word2Vec	0.659	0.541	0.594
Oceanic CFBF	0.705	0.648	0.676
Oceanic Word2Vec + Oceanic CFBF	0.705	0.604	0.650
Flair + GloVe (General-purpose) + Oceanic Word2Vec + Oceanic CFBF	0.713	0.649	0.679

^a In this task, a true positive result is one where the algorithm correctly identifies the named entity and assigns the correct class.

843 challenges facing the ocean science community is to integrate this vast amount
 844 of data in a way that will facilitate its translation into improved understanding
 845 of oceanic processes. Addressing this challenge relies strongly on the implemen-
 846 tation of AI technologies, which now, in the era of Big Data, are ubiquitously
 847 applied across scientific domains and disciplines.

848 In this paper, we have deconstructed the process of oceanic science DI and
 849 pointed to the key missing tools and underutilized information sources currently
 850 limiting its automation. We have focused on semantic AI technologies aiding the
 851 matching and mapping phases of the DI process, limiting our discussion of data
 852 fusion and data cleansing techniques, which we intend to address in future work.
 853 The potential of implementing AI technologies to advance oceanic research calls
 854 for close collaboration between ocean and data scientists. Importantly, such col-
 855 laboration should promote the formation of dedicated infrastructures to support
 856 AI efforts in ocean science, focusing on several activities that address major limi-
 857 tations in the current state of ocean data integration (Table 3):

- 858 • Develop AI-based tools for assisting ocean scientists in aligning their
 859 schema with existing ontologies when organizing their measurements in
 860 datasets.
- 861 • Extend and refine conceptual coverage of – and conceptual alignment be-
 862 tween – existing ontologies, such that they are more compatible with the
 863 diverse and multidisciplinary nature of ocean science.
- 864 • Create ocean-science-specific schema matching and mapping benchmarks
 865 to accelerate the development of matching and mapping tools utilizing se-

- 866 mantics encoded in existing vocabularies and ontologies.
- 867 • Similarly support the development of ocean-science-specific entity resolu-
- 868 tion tools by creating annotated datasets and benchmarks on both the dataset
- 869 and data point level.
- 870 • Annotate datasets, and develop tools and benchmarks for the extraction and
- 871 categorization of data quality and preprocessing descriptions from scientific
- 872 text.
- 873 • Create large-scale word embeddings trained upon ocean science literature
- 874 to accelerate the development of AI-based information extraction, entity
- 875 resolution, and matching tools.

876 Formation of improved AI integration infrastructure based on these suggested

877 activities will contribute importantly to our ability to share, explore, and inter-

878 pret the vast amount of available oceanic data, thus substantially advancing ocean

879 research.

880 **6. Author contributions**

881 T.S. led sections 2 and 3, Y.L. led sections 1 and 5, and K.B. led section 4. The

882 deconstruction of a DI process in ocean science as portrayed in Figure 4 was

883 performed jointly by T.S. and Y.L.

884 **7. Conflicts of interest**

885 None identified.

886 **References**

- 887 Abedjan Z, Chu X, Deng D, Fernandez RC, Ilyas IF, Ouzzani M, Papotti P,
- 888 Stonebraker M, Tang N. 2016. Detecting Data Errors: Where are we and
- 889 what needs to be done? *PVLDB* **9**(12): 993–1004. doi:10.14778/2994509-
- 890 2994518.
- 891 Akbik A, Blythe D, Vollgraf R. 2018. Contextual String Embeddings for Sequence
- 892 Labeling, in *Proceedings of the 27th International Conference on Computa-*
- 893 *tional Linguistics, COLING 2018, Santa Fe, New Mexico, USA, August 20-*
- 894 *26, 2018*, pp. 1638–1649. [https://www.aclweb.org/anthology/](https://www.aclweb.org/anthology/C18-1139/)
- 895 [C18-1139/](https://www.aclweb.org/anthology/C18-1139/).
- 896 Alexe B, ten Cate B, Kolaitis PG, Tan WC. 2011. EIRENE: Interactive Design and
- 897 Refinement of Schema Mappings via Data Examples. *PVLDB* **4**(12): 1414–
- 898 1417. <http://www.vldb.org/pvldb/vol4/p1414-alexe.pdf>.

- 899 Anaconda. 2020. Anaconda Distribution. Retrieved Jan. 22nd, 2020. [https://](https://www.anaconda.com/distribution/)
900 www.anaconda.com/distribution/.
- 901 Ashish N. 2005. Semantic-Web Technology: Applications at NASA, in Kalfoglou
902 Y, Schorlemmer M, Sheth A, Staab S, Uschold M, eds., *Semantic Interop-*
903 *erability and Integration*. Dagstuhl, Germany: Internationales Begegnungs-
904 und Forschungszentrum für Informatik (IBFI), Schloss Dagstuhl, Ger-
905 many. (Dagstuhl Seminar Proceedings 04391). ISSN 1862-4405. [http:](http://drops.dagstuhl.de/opus/volltexte/2005/32)
906 [//drops.dagstuhl.de/opus/volltexte/2005/32](http://drops.dagstuhl.de/opus/volltexte/2005/32).
- 907 Assmy P, Henjes J, Klaas C, Smetacek V. 2007. Mechanisms determining species
908 dominance in a phytoplankton bloom induced by the iron fertilization exper-
909 iment EisenEx in the Southern Ocean. *Deep-Sea Res Part I-Oceanogr Res*
910 *Pap* **54**(3): 340–362.
- 911 Auer S, Bizer C, Kobilarov G, Lehmann J, Cyganiak R, Ives ZG. 2007. DBpedia:
912 A Nucleus for a Web of Open Data, in Aberer K, Choi K, Noy NF, Alle-
913 mang D, Lee K, Nixon LJB, Golbeck J, Mika P, Maynard D, Mizoguchi R,
914 Schreiber G, Cudré-Mauroux P, eds., *The Semantic Web, 6th International*
915 *Semantic Web Conference, 2nd Asian Semantic Web Conference, ISWC 2007*
916 *+ ASWC 2007, Busan, Korea, November 11-15, 2007*, vol. 4825, pp. 722–
917 735. Springer. doi:10.1007/978-3-540-76298-0_52.
- 918 Bar K. 2020a. Oceanic NER Project. Retrieved Jan. 22nd, 2020. doi:10.17605/
919 OSF.IO/MY2NK.
- 920 Bar K. 2020b. Oceanic Data Description Extraction Project. Retrieved Jan. 22nd,
921 2020. doi:10.17605/OSF.IO/8VAFS.
- 922 Bellahsene Z, Bonifati A, Rahm E, eds. 2011. *Schema Matching and Map-*
923 *ping*. Berlin, Heidelberg: Springer. (Data-Centric Systems and Applications).
924 ISBN 978-3-642-16517-7. doi:10.1007/978-3-642-16518-4.
- 925 Berg JL. 1976. Data base directions: the next steps. *ACM SIGMOD Record* **8**(2):
926 3–4.
- 927 Berners-Lee T, Hendler J. 2001. Publishing on the semantic web. *Nature*
928 **410**(6832): 1023–1024. doi:10.1038/35074206.
- 929 Biological and Chemical Oceanography Data Management Office. 2020. Intro-
930 duction to BCO-DMO. Retrieved Jan. 3rd, 2020. [https://www.bco-](https://www.bco-dmo.org/)
931 [dmo.org/](https://www.bco-dmo.org/).
- 932 Bolukbasi T, Chang KW, Zou JY, Saligrama V, Kalai AT. 2016. Man is
933 to Computer Programmer as Woman is to Homemaker? Debias-
934 ing Word Embeddings, in Lee DD, Sugiyama M, von Luxburg U,
935 Guyon I, Garnett R, eds., *Advances in Neural Information Process-*
936 *ing Systems 29: Annual Conference on Neural Information Processing*

- 937 *Systems 2016, December 5-10, 2016, Barcelona, Spain*, pp. 4349–
938 4357. [http://papers.nips.cc/paper/6228-man-is-to-](http://papers.nips.cc/paper/6228-man-is-to-computer-programmer-as-woman-is-to-homemaker-debiasing-word-embeddings)
939 [computer-programmer-as-woman-is-to-homemaker-](http://papers.nips.cc/paper/6228-man-is-to-computer-programmer-as-woman-is-to-homemaker-debiasing-word-embeddings)
940 [debiasing-word-embeddings](http://papers.nips.cc/paper/6228-man-is-to-computer-programmer-as-woman-is-to-homemaker-debiasing-word-embeddings).
- 941 British Oceanographic Data Centre. 2020. European Directory of Marine Environ-
942 mental Data. Retrieved Jan. 3rd, 2020. [https://edmed.seadatanet.](https://edmed.seadatanet.org/)
943 [org/](https://edmed.seadatanet.org/).
- 944 Chen Z, Jia H, Heflin J, Davison BD. 2018. Generating Schema Labels Through
945 Dataset Content Analysis, in *Companion Proc. of the The Web Conf. 2018*,
946 pp. 1515–1522. Republic and Canton of Geneva, Switzerland: International
947 World Wide Web Conferences Steering Committee. (WWW '18). doi:10.
948 1145/3184558.3191601.
- 949 Claramunt C, Ray C, Salmon L, Camossi E, Hadzagic M, Jousset A-L, An-
950 drienko G, Andrienko N, Theodoridis Y, Vouros G. 2017. Maritime data inte-
951 gration and analysis: recent progress and research challenges. in Markl V, Or-
952 lando S, Mitschang B, Andritsos P, Sattler K-U, Breß S, eds., *Proceedings of*
953 *the 20th International Conference on Extending Database Technology, EDBT*
954 *2017, Venice, Italy, March 21-24, 2017*, pp. 192–197. OpenProceedings.org.
955 doi:10.5441/002/edbt.2017.18.
- 956 Clark CA, Divvala S. 2015. Looking Beyond Text: Extracting Figures, Tables,
957 and Captions from Computer Science Papers, in *Workshops at the Twenty-*
958 *Ninth AAAI Conference on Artificial Intelligence, Austin, TX, USA, January*
959 *25 - 26, 2015*, vol. 53, pp. 599–605. [https://www.aaai.org/ocs/](https://www.aaai.org/ocs/index.php/WS/AAAIW15/paper/viewPaper/10092)
960 [index.php/WS/AAAIW15/paper/viewPaper/10092](https://www.aaai.org/ocs/index.php/WS/AAAIW15/paper/viewPaper/10092).
- 961 Claus S, De Hauwere N, Vanhoorne B, Deckers P, Souza Dias F, Hernan-
962 dez F, Mees J. 2014. Marine regions: towards a global standard for geo-
963 referenced marine names and boundaries. *Mar Geod* **37**(2): 99–125. doi:
964 10.1080/01490419.2014.902881.
- 965 Data Observation Network for Earth. 2020. The Patience of the Data
966 Hunter. Retrieved Jan. 3rd, 2020. [https://www.dataone.org/](https://www.dataone.org/data-stories/patience-data-hunter)
967 [data-stories/patience-data-hunter](https://www.dataone.org/data-stories/patience-data-hunter).
- 968 De Uña D, Rümmele N, Gange G, Schachte P, Stuckey PJ. 2018. Machine Learn-
969 ing and Constraint Programming for Relational-to-Ontology Schema Map-
970 ping, in *Proceedings of the 27th International Joint Conference on Artificial*
971 *Intelligence*, pp. 1277–1283. AAAI Press. (IJCAI'18). doi:10.5555/3304415.
972 3304597.
- 973 Devlin J, Chang M-W, Lee K, Toutanova K. 2019. BERT: Pre-training of Deep
974 Bidirectional Transformers for Language Understanding, in Burstein J, Do-

- 975 ran C, Solorio T, eds., *Proceedings of the 2019 Conference of the North*
976 *American Chapter of the Association for Computational Linguistics: Hu-*
977 *man Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA,*
978 *June 2-7, 2019, Volume 1 (Long and Short Papers)*, pp. 4171–4186. As-
979 sociation for Computational Linguistics. [https://www.aclweb.org/](https://www.aclweb.org/anthology/N19-1423/)
980 [anthology/N19-1423/](https://www.aclweb.org/anthology/N19-1423/).
- 981 Do HH, Rahm E. 2002. COMA - A System for Flexible Combination of Schema
982 Matching Approaches, in *Proceedings of 28th International Conference on*
983 *Very Large Data Bases, VLDB 2002, Hong Kong, August 20-23, 2002*, pp.
984 610–621. Morgan Kaufmann. doi:10.1016/B978-155860869-6/50060-3.
- 985 Doan A, Domingos PM, Levy AY. 2000. Learning Source Description for
986 Data Integration, in Suciu D, Vossen G, eds., *Proceedings of the*
987 *Third International Workshop on the Web and Databases, WebDB 2000,*
988 *Adam’s Mark Hotel, Dallas, Texas, USA, May 18-19, 2000, in con-*
989 *junction with ACM PODS/SIGMOD 2000. Informal proceedings*, pp.
990 81–86. [http://citeseerx.ist.psu.edu/viewdoc/summary?](http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.134.3677)
991 [doi=10.1.1.134.3677](http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.134.3677).
- 992 Doan A, Madhavan J, Domingos PM, Halevy AY. 2002. Learning to map be-
993 tween ontologies on the semantic web, in Lassner D, Roure DD, Iyengar A,
994 eds., *Proceedings of the Eleventh International World Wide Web Conference,*
995 *WWW 2002, May 7-11, 2002, Honolulu, Hawaii, USA*, pp. 662–673. ACM.
996 doi:10.1145/511446.511532.
- 997 Dong XL, Rekatsinas T. 2018. Data Integration and Machine Learning: A Nat-
998 ural Synergy, in Das G, Jermaine CM, Bernstein PA, eds., *Proceedings of*
999 *the 2018 International Conference on Management of Data, SIGMOD Con-*
1000 *ference 2018, Houston, TX, USA, June 10-15, 2018*, pp. 1645–1650. ACM.
1001 doi:10.1145/3183713.3197387.
- 1002 Dong XL, Srivastava D. 2015. *Big Data Integration*. Morgan & Claypool Publish-
1003 ers. doi:10.2200/S00578ED1V01Y201404DTM040.
- 1004 Durden JM, Luo JY, Alexander H, Flanagan AM, Grossmann L. 2017. Integrat-
1005 ing “Big Data” into aquatic ecology: challenges and opportunities. *Limnol*
1006 *Oceanogr Bull* **26**(4): 101–108. doi:10.1002/lob.10213.
- 1007 Ebraheem M, Thirumuruganathan S, Joty SR, Ouzzani M, Tang N. 2018. Dis-
1008 tributed representations of tuples for entity resolution. *PVLDB* **11**(11):
1009 1454–1467. doi:10.14778/3236187.3236198. [http://www.vldb.org/](http://www.vldb.org/pvldb/vol11/p1454-ebraheem.pdf)
1010 [pvldb/vol11/p1454-ebraheem.pdf](http://www.vldb.org/pvldb/vol11/p1454-ebraheem.pdf).
- 1011 Ekaputra FJ, Sabou M, Serral E, Kiesling E, Biffi S. 2017. Ontology-based data
1012 integration in multi-disciplinary engineering environments: A Review. *Open*

- 1013 *Journal of Information Systems (OJIS)* **4**(1): 1–26.
- 1014 Eriksen CC, Osse TJ, Light RD, Wen T, Lehman TW, Sabin PL, Ballard JW,
1015 Chiodi AM. 2001. Seaglider: A long-range autonomous underwater vehicle
1016 for oceanographic research. *IEEE J Ocean Eng* **26**(4): 424–436.
- 1017 European Commission. 2020. Copernicus, the European Earth Observation and
1018 Monitoring Programme. Retrieved Jan. 1st, 2020. [http://copernicus.
1019 eu/](http://copernicus.eu/).
- 1020 Fernandez RC, Mansour E, Qahtan AA, Elmagarmid AK, Ilyas IF, Madden S,
1021 Ouzzani M, Stonebraker M, Tang N. 2018. Seeping Semantics: Linking
1022 Datasets Using Word Embeddings for Data Discovery. in *34th IEEE Interna-
1023 tional Conference on Data Engineering, ICDE 2018, Paris, France, April 16-
1024 19, 2018*, pp. 989–1000. IEEE Computer Society. doi:10.1109/ICDE.2018.
1025 00093. <https://doi.org/10.1109/ICDE.2018.00093>.
- 1026 Field CB, Behrenfeld MJ, Randerson JT, Falkowski P. 1998. Primary produc-
1027 tion of the biosphere: integrating terrestrial and oceanic components. *Science*
1028 **281**(5374): 237–240. ISSN 0036-8075. doi:10.1126/science.281.5374.237.
- 1029 Freund Y, Mason L. 1999. The Alternating Decision Tree Learning Algorithm, in
1030 *Proceedings of the Sixteenth International Conference on Machine Learning*,
1031 pp. 124–133. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.
1032 (ICML '99). ISBN 1558606122.
- 1033 Froese R, Pauly D. 2020. FishBase. Retrieved Jan. 8th, 2020. [https://www.
1034 fishbase.ca](https://www.fishbase.ca).
- 1035 Gal A, Sagi T. 2010. Tuning the ensemble selection process of schema matchers.
1036 *Inf Syst* **35**(8): 845–859. doi:10.1016/j.is.2010.04.003.
- 1037 Gal A. 2011. *Uncertain Schema Matching*. Morgan & Claypool
1038 Publishers. (Synthesis Lectures on Data Management). doi:
1039 10.2200/S00337ED1V01Y201102DTM013.
- 1040 Gangemi A. 2005. Ontology Design Patterns for Semantic Web Content, in Gil Y,
1041 Motta E, Benjamins VR, Musen MA, eds., *The Semantic Web – ISWC 2005*,
1042 pp. 262–276. Berlin, Heidelberg: Springer Berlin Heidelberg. doi:10.1007/
1043 11574620_21.
- 1044 Goodhue DL, Wybo MD, Kirsch LJ. 1992. The impact of data integration on
1045 the costs and benefits of information systems. *MIS Q* **16**(3): 293–311.
1046 [http://misq.org/the-impact-of-data-integration-on-
1047 the-costs-and-benefits-of-information-systems.html](http://misq.org/the-impact-of-data-integration-on-the-costs-and-benefits-of-information-systems.html).
- 1048 Gregory K, Groth P, Cousijn H, Scharnhorst A, Wyatt S. 2019. Searching data:
1049 a review of observational data retrieval practices in selected disciplines. *J
1050 Assoc Inf Sci Tech* **70**(5): 419–432. doi:10.1002/asi.24165.

- 1051 Gruber TR. 1995. Toward principles for the design of ontologies used for knowl-
1052 edge sharing? *Int J Hum-Comput Stud* **43**(5-6): 907–928. doi:10.1006/ijhc.
1053 1995.1081.
- 1054 Gubanov MN, Stonebraker M, Bruckner D. 2014. Text and structured data fusion
1055 in data tamer at scale, in Cruz IF, Ferrari E, Tao Y, Bertino E, Trajcevski
1056 G, eds., *IEEE 30th International Conference on Data Engineering, Chicago,*
1057 *ICDE 2014, IL, USA, March 31 - April 4, 2014*, pp. 1258–1261. IEEE Com-
1058 puter Society. doi:10.1109/ICDE.2014.6816755.
- 1059 Guiry MD, Guiry GM. 2020. AlgaeBase. World-wide electronic publication, Na-
1060 tional University of Ireland, Galway. Searched on Jan. 8th, 2020. [https:](https://www.algaebase.org)
1061 [//www.algaebase.org](https://www.algaebase.org).
- 1062 Gupta M, Gao J, Aggarwal CC, Han J. 2014. Outlier detection for temporal data:
1063 A survey. *IEEE Trans Knowl Data Eng* **26**(9): 2250–2267. ISSN 2326-3865.
1064 doi:10.1109/TKDE.2013.184.
- 1065 Halevy AY, Norvig P, Pereira F. 2009. The unreasonable effectiveness of data.
1066 *IEEE Intell Syst* **24**(2): 8–12. doi:10.1109/MIS.2009.36.
- 1067 Halevy AY, Rajaraman A, Ordille JJ. 2006. Data Integration: The Teenage Years,
1068 in Dayal U, Whang K, Lomet DB, Alonso G, Lohman GM, Kersten ML, Cha
1069 SK, Kim Y, eds., *Proceedings of the 32nd International Conference on Very*
1070 *Large Data Bases, Seoul, Korea, September 12-15, 2006*, pp. 9–16. ACM.
1071 <http://dl.acm.org/citation.cfm?id=1164130>.
- 1072 Hammer M, McLeod D. 1979. On Database Management System Architec-
1073 ture. Defense Technical Information Center. [http://www.dtic.mil/](http://www.dtic.mil/docs/citations/ADA076417)
1074 [docs/citations/ADA076417](http://www.dtic.mil/docs/citations/ADA076417).
- 1075 Hartigan JA, Wong MA. 1979. Algorithm AS 136: A k-means clustering algo-
1076 rithm. *J R Stat Soc Ser C-Appl Stat* **28**(1): 100–108.
- 1077 He B, Chang KC-C. 2006. Automatic complex schema matching across Web
1078 query interfaces: A correlation mining approach. *ACM Trans Database Syst*
1079 **31**(1): 346–395. doi:10.1145/1132863.1132872.
- 1080 Hinton G, Deng L, Yu D, Dahl GE, Mohamed A, Jaitly N, Senior A, Vanhoucke V,
1081 Nguyen P, Sainath TN, Kingsbury B. 2012. Deep neural networks for acous-
1082 tic modeling in speech recognition: the shared views of four research groups.
1083 *IEEE Signal Process Mag* **29**(6): 82–97. doi:10.1109/MSP.2012.2205597.
- 1084 Hogan A, Zimmermann A, Umbrich J, Polleres A, Decker S. 2012. Scalable and
1085 distributed methods for entity matching, consolidation and disambiguation
1086 over linked data corpora. *J Web Semant* **10**: 76–110.
- 1087 IPCC. 2014. *Climate Change 2014: Synthesis Report. Contribution of Working*
1088 *Groups I, II and III to the Fifth Assessment Report of the Intergovernmental*

- 1089 *Panel on Climate Change*. [Core Writing Team, Pachauri RK and Meyer LA
1090 (eds.)]. IPCC, Geneva, Switzerland. 151 pp. doi:10013/epic.45156.d001.
- 1091 Jiang R, Banchs RE, Li H. 2016. Evaluating and Combining Name Entity Recog-
1092 nition Systems, in *Proceedings of the Sixth Named Entity Workshop*, pp.
1093 21–27. Berlin, Germany: Association for Computational Linguistics. doi:
1094 10.18653/v1/W16-2703.
- 1095 Jordan MI, Mitchell TM. 2015. Machine learning: Trends, perspectives, and
1096 prospects. *Science* **349**(6245): 255–60. doi:10.1126/science.aaa8415.
- 1097 Kaplan A, Haenlein M. 2019. Siri, Siri, in my hand: Who’s the fairest in the land?
1098 On the interpretations, illustrations, and implications of artificial intelligence.
1099 *Bus Horiz* **62**(1): 15 – 25. doi:10.1016/j.bushor.2018.08.004.
- 1100 Kaplan DM, Lekien F. 2007. Spatial interpolation and filtering of surface cur-
1101 rent data based on open-boundary modal analysis. *J Geophys Res* **112**(C12):
1102 C12007. doi:10.1029/2006JC003984.
- 1103 Kenig B, Gal A. 2013. MFIBlocks: An effective blocking algorithm for entity
1104 resolution. *Inf Syst* **38**(6): 908–926. doi:10.1016/j.is.2012.11.008.
- 1105 Krisnadhi AA, Hu Y, Janowicz K, Hitzler P, Arko RA, Carbotte S, Chandler C,
1106 Cheatham M, Fils D, Finin T, Ji P, Jones MB, Karima N, Lehnert KA, Mickle
1107 A, Narock T, O’Brien M, Raymond L, Shepherd A, Schildhauer M, Wiebe
1108 P. 2015. The GeoLink Framework for Pattern-based Linked Data Integration,
1109 in Villata S, Pan JZ, Dragoni M, eds., *Proceedings of the ISWC 2015 Posters
1110 & Demonstrations Track co-located with the 14th International Semantic
1111 Web Conference (ISWC-2015), Bethlehem, PA, USA, October 11, 2015,
1112 CEUR Workshop Proceedings*, vol. 1486. CEUR-WS.org. (CEUR Workshop
1113 Proceedings, vol. 1486). [http://ceur-ws.org/Vol-1486/paper_](http://ceur-ws.org/Vol-1486/paper_99.pdf)
1114 [99.pdf](http://ceur-ws.org/Vol-1486/paper_99.pdf).
- 1115 Krizhevsky A, Sutskever I, Hinton GE. 2017. ImageNet classification with deep
1116 convolutional neural networks. *Commun ACM* **60**(6): 84–90. doi:10.1145/
1117 3065386.
- 1118 Lammey R. 2015. CrossRef Text and Data Mining Services. *Science Editing* **2**:
1119 22–27. doi:10.6087/kcse.32.
- 1120 Leadbetter A, Hamre T, Lowry R, Lassoued Y, Dunne D. 2010. Ontologies and
1121 ontology extension for marine environmental information systems, in Arne
1122 J Berre DR, Maue P, eds., *Proceedings of the Workshop Environmental In-
1123 formation Systems and Services-Infrastructures and Platforms,(envip’2010),
1124 Bonn, Germany*, vol. 34(25), pp. 12–14. [http://ceur-ws.org/Vol-](http://ceur-ws.org/Vol-679/paper11.pdf)
1125 [679/paper11.pdf](http://ceur-ws.org/Vol-679/paper11.pdf).
- 1126 Leblanc K, Arístegui J, Armand L, Assmy P, Beker B, Bode A, Breton E, Cor-

- 1127 net V, Gibson J, Gosselin MP, Kopczynska E, Marshall H, Peloquin J, Pi-
1128 ontkovski S, Poulton AJ, Quéguiner B, Schiebel R, Shipe R, Stefels J, van
1129 Leeuwe MA, Varela M, Widdicombe C, Yallop M. 2012. A global diatom
1130 database – abundance, biovolume and biomass in the world ocean. *Earth*
1131 *Syst Sci Data* **4**: 149–165. doi:10.5194/essd-4-149-2012.
- 1132 Lehahn Y, d’Ovidio F, Koren I. 2018. A satellite-based lagrangian view on phyto-
1133 plankton dynamics. *Annu Rev Mar Sci* **10**: 99–119.
- 1134 Lehahn Y, Ingle KN, Golberg A. 2016. Global potential of offshore and shallow
1135 waters macroalgal biorefineries to provide for food, chemicals and energy:
1136 feasibility and sustainability. *Algal Res* **17**: 150–160.
- 1137 Lesiv M, Moltchanova E, Schepaschenko D, See L, Shvidenko A, Comber A,
1138 Fritz S. 2016. Comparison of data fusion methods using crowdsourced data
1139 in creating a hybrid forest cover map. *Remote Sens* **8**(3): 261.
- 1140 Levy O, Goldberg Y, Dagan I. 2015. Improving distributional similarity with
1141 lessons learned from word embeddings. *TACL* **3**: 211–225. doi:10.1162/
1142 tacl.a_00134.
- 1143 Lumpkin R, Özgökmen T, Centurioni L. 2017. Advances in the application of
1144 surface drifters. *Annu Rev Mar Sci* **9**: 59–81.
- 1145 Luo Y-W, Doney SC, Anderson LA, Benavides M, Berman-Frank I, Bode A,
1146 Bonnet S, Boström KH, Böttjer D, Capone DG, Carpenter EJ, Chen YL,
1147 Church MJ, Dore JE, Falcón LI, Fernández A, Foster RA, Furuya K, Gómez
1148 F, Gundersen K, Hynes AM, Karl DM, Kitajima S, Langlois RJ, LaRoche
1149 J, Letelier RM, Marañón E, McGillicuddy DJ, Moisander PH, Moore CM,
1150 Mouriño-Carballido B, Mulholland MR, Needoba JA, Orcutt KM, Poulton
1151 AJ, Rahav E, Raimbault P, Rees AP, Riemann L, Shiozaki T, Subrama-
1152 niam A, Tyrrell T, Turk-Kubo KA, Varela M, Villareal TA, Webb EA, White
1153 AE, Wu J, Zehr JP. 2012. Database of diazotrophs in global ocean: abun-
1154 dances, biomass and nitrogen fixation rates. *Earth Syst Sci Data* **4**: 47–73.
1155 doi:10.5194/essd-4-47-2012.
- 1156 Madhavan J, Bernstein PA, Doan A, Halevy AY. 2005. Corpus-based Schema
1157 Matching, in Aberer K, Franklin MJ, Nishio S, eds., *Proceedings of the 21st*
1158 *International Conference on Data Engineering, ICDE 2005, 5-8 April 2005,*
1159 *Tokyo, Japan*, pp. 57–68. IEEE Computer Society. doi:10.1109/ICDE.2005.
1160 39.
- 1161 Martinez-Rodriguez JL, Hogan A, Lopez-Arevalo I. 2020. Information extraction
1162 meets the semantic web: a survey. *Semant Web* **11**(2): 255–335. doi:10.3233/
1163 SW-180333.
- 1164 Mayhew S, Roth D. 2018. TALEN: Tool for Annotation of Low-resource ENTities,

- 1165 in *Proceedings of the 56th Annual Meeting of the Association for Computa-*
1166 *tional Linguistics-System Demonstrations, Melbourne, Australia, July 15 -*
1167 *20, 2018*, pp. 80–86.
- 1168 Mikolov T, Yih W, Zweig G. 2013. Linguistic Regularities in Continuous Space
1169 Word Representations, in Vanderwende L, III HD, Kirchhoff K, eds., *Hu-*
1170 *man Language Technologies: Conference of the North American Chapter of*
1171 *the Association of Computational Linguistics, Proceedings, June 9-14, 2013,*
1172 *Westin Peachtree Plaza Hotel, Atlanta, Georgia, USA*, pp. 746–751. The As-
1173 sociation for Computational Linguistics. [https://www.aclweb.org/](https://www.aclweb.org/anthology/N13-1090/)
1174 [anthology/N13-1090/](https://www.aclweb.org/anthology/N13-1090/).
- 1175 Narock T, Arko RA, Carbotte S, Krisnadhi A, Hitzler P, Cheatham M, Shepherd A,
1176 Chandler C, Raymond L, Wiebe P, Finin TW. 2014. The OceanLink project,
1177 in Lin JJ, Pei J, Hu X, Chang W, Nambiar R, Aggarwal CC, Cercone N,
1178 Honavar VG, Huan J, Mobasher B, Pyne S, eds., *2014 IEEE International*
1179 *Conference on Big Data, Big Data 2014, Washington, DC, USA, October 27-*
1180 *30, 2014*, pp. 14–21. IEEE Computer Society. doi:10.1109/BigData.2014.
1181 7004347.
- 1182 National Oceanic and Atmospheric Administration. 2020a. Big Data Project.
1183 Retrieved Jan. 3rd, 2020. [https://www.noaa.gov/big-data-](https://www.noaa.gov/big-data-project)
1184 [project](https://www.noaa.gov/big-data-project).
- 1185 National Oceanic and Atmospheric Administration. 2020b. National Centers
1186 for Environmental Information. Retrieved Jan. 1st, 2020. [https://www.](https://www.ncei.noaa.gov/)
1187 [ncei.noaa.gov/](https://www.ncei.noaa.gov/).
- 1188 Nickel M, Murphy K, Tresp V, Gabrilovich E. 2016. A review of relational
1189 machine learning for knowledge graphs. *Proc IEEE* **104**(1): 11–33. doi:
1190 10.1109/JPROC.2015.2483592.
- 1191 O’Brien CJ, Pelloquin JA, Vogt M, Heinle M, Gruber N, Ajani P, Andrleit H,
1192 Arístegui J, Beaufort L, Estrada M, Karentz D, Kopczyńska E, Lee R, Poul-
1193 ton AJ, Pritchard T, Widdicombe C. 2013. Global marine plankton functional
1194 type biomass distributions: Coccolithophores. *Earth Syst Sci Data* **5**(2): 259–
1195 276. doi:10.5194/essd-5-259-2013.
- 1196 O’Hare K, Jurek-Loughrey A, de Campos C. 2019. A Review of Unsupervised
1197 and Semi-supervised Blocking Methods for Record Linkage, in Deepak P,
1198 Jurek-Loughrey A, eds., *Linking and Mining Heterogeneous and Multi-view*
1199 *Data*. Cham: Springer International Publishing: pp. 79–105. ISBN 978-3-
1200 030-01872-6. doi:10.1007/978-3-030-01872-6_4.
- 1201 PANGAEA. 2020. PANGAEA, Data Publisher for Earth and Environmental Science.
1202 Retrieved Jan. 1st, 2020. <https://pangaea.de/>.

- 1203 Papadakis G, Svirsky J, Gal A, Palpanas T. 2016. Comparative analysis of ap-
1204 proximate blocking techniques for entity resolution. *PVLDB* **9**(9): 684–695.
1205 doi:10.14778/2947618.2947624.
- 1206 Pennington J, Socher R, Manning CD. 2014. Glove: Global Vectors for Word
1207 Representation, in Moschitti A, Pang B, Daelemans W, eds., *Proceedings of*
1208 *the 2014 Conference on Empirical Methods in Natural Language Processing,*
1209 *EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a*
1210 *Special Interest Group of the ACL*, pp. 1532–1543. ACL. <https://www.aclweb.org/anthology/D14-1162/>.
- 1212 Pennington J, Socher R, Manning CD. 2020. GloVe: Global Vectors for Word
1213 Representation. Retrieved Jan. 22nd, 2020. <https://nlp.stanford.edu/projects/glove/>.
- 1215 Prud'hommeaux E, Seaborne A. 2008. SPARQL Query Language for RDF. W3C.
1216 <http://www.w3.org/TR/rdf-sparql-query/>.
- 1217 Řehůřek R, Sojka P. 2010. Software Framework for Topic Modelling with Large
1218 Corpora, in *Proceedings of the LREC 2010 Workshop on New Challenges for*
1219 *NLP Frameworks*, pp. 46–50. Valletta, Malta: ELRA. <http://is.muni.cz/publication/884893/en>.
- 1221 Roemmich D, Johnson GC, Riser S, Davis R, Gilson J, Owens WB, Garzoli SL,
1222 Schmid C, Ignaszewski M. 2009. The Argo Program: Observing the global
1223 ocean with profiling floats. *Oceanogr* **22**: 34–43.
- 1224 Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpa-
1225 thy A, Khosla A, Bernstein MS, Berg AC, Li F. 2015. ImageNet large
1226 scale visual recognition challenge. *Int J Comput Vis* **115**(3): 211–252. doi:
1227 10.1007/s11263-015-0816-y.
- 1228 Sagi T, Gal A. 2013. Schema matching prediction with applications to data source
1229 discovery and dynamic ensembling. *VLDB J* **22**(5): 689–710. doi:10.1007/
1230 s00778-013-0325-y.
- 1231 Sagi T, Gal A, Barkol O, Bergman R, Avram A. 2017. Multi-source uncertain
1232 entity resolution: transforming holocaust victim reports into people. *Inf Syst*
1233 **65**: 124–136. doi:10.1016/j.is.2016.12.003.
- 1234 Sang EFTK, De Meulder F. 2003. Introduction to the CoNLL-2003 Shared Task:
1235 Language-Independent Named Entity Recognition, in Daelemans W, Os-
1236 borne M, eds., *Proceedings of the Seventh Conference on Natural Language*
1237 *Learning, CoNLL 2003, Held in cooperation with HLT-NAACL 2003, Ed-*
1238 *monton, Canada, May 31 - June 1, 2003*, pp. 142–147. ACL. <https://www.aclweb.org/anthology/W03-0419/>.
- 1240 Schmidt EM, Kim YE. 2011. Learning emotion-based acoustic features with deep

- 1241 belief networks, in *IEEE Workshop on Applications of Signal Processing to*
1242 *Audio and Acoustics, WASPAA 2011, New Paltz, NY, USA, October 16-19,*
1243 *2011*, pp. 65–68. IEEE. doi:10.1109/ASPAA.2011.6082328.
- 1244 Semina GI, Mikaelyan AS. 1994. (Table 1) Hydrological, hydrooptical, and hy-
1245 drochemical characteristics of seawater at 7 stations in the Northwest Pacific.
1246 PANGAEA. doi:10.1594/PANGAEA.759517. In supplement to: Semina, GI;
1247 Mikaelyan, AS (1994): Phytoplankton of various size groups from the North-
1248 west Pacific Ocean during summer. *Oceanology*, 33(5), 618-624.
- 1249 Shvaiko P, Euzenat J. 2013. Ontology matching: state of the art and future chal-
1250 lenges. *IEEE Trans Knowl Data Eng* **25**(1): 158–176. doi:10.1109/TKDE.
1251 2011.253.
- 1252 Silver D, Huang A, Maddison CJ, Guez A, Sifre L, van den Driessche G,
1253 Schrittwieser J, Antonoglou I, Panneershelvam V, Lanctot M, Dieleman S,
1254 Grewe D, Nham J, Kalchbrenner N, Sutskever I, Lillicrap TP, Leach M,
1255 Kavukcuoglu K, Graepel T, Hassabis D. 2016. Mastering the game of Go
1256 with deep neural networks and tree search. *Nature* **529**(7587): 484–489. doi:
1257 10.1038/nature16961.
- 1258 Sorrentino S, Bergamaschi S, Gawinecki M, Po L. 2010. Schema label normaliza-
1259 tion for improving schema matching. *Data Knowl Eng* **69**(12): 1254–1273.
1260 doi:10.1016/j.datak.2010.10.004.
- 1261 Stocker TF, Qin D, Plattner GK, Alexander LV, Allen SK, Bindoff NL, Bréon FM,
1262 Church JA, Cubasch U, Emori S, Forster P, Friedlingstein P, Gillett N, Gre-
1263 gory JM, Hartmann DL, Jansen E, Kirtman B, Knutti R, Krishna Kumar K,
1264 Lemke P, Marotzke J, Masson-Delmotte V, Meehl GA, Mokhov II, Piao S,
1265 Ramaswamy V, Randall D, Rhein M, Rojas M, Sabine C, Shindell D, Talley
1266 LD, Vaughan DG, Xie SP. 2013. Technical Summary, in Stocker T, Qin D,
1267 Plattner GK, Tignor M, Allen S, Boschung J, Nauels A, Xia Y, Bex V, Midg-
1268 ley P, eds., *Climate Change 2013: The Physical Science Basis. Contribution*
1269 *of Working Group I to the Fifth Assessment Report of the Inter-governmental*
1270 *Panel on Climate Change*. Cambridge, United Kingdom and New York, NY,
1271 USA: Cambridge University Press.
- 1272 Tzitzikas Y, Allocca C, Bekiari C, Marketakis Y, Fafalios P, Doerr M, Minadakis
1273 N, Patkos T, Candela L. 2013. Integrating Heterogeneous and Distributed In-
1274 formation about Marine Species through a Top Level Ontology, in Garoufal-
1275 lou E, Greenberg J, eds., *Metadata and Semantics Research - 7th Research*
1276 *Conference, MTSR 2013, Thessaloniki, Greece, November 19-22, 2013. Pro-*
1277 *ceedings, Communications in Computer and Information Science*, vol. 390,
1278 pp. 289–301. Springer. (Communications in Computer and Information Sci-

- 1279 ence, vol. 390). doi:10.1007/978-3-319-03437-9_29.
- 1280 UNIDATA. 2019. Network Common Data Form (NetCDF). Retrieved Jan. 3rd,
1281 2020. <https://www.unidata.ucar.edu/software/netcdf/>.
- 1282 Uschold M. 1998. Knowledge level modelling: concepts and terminology. *Knowl*
1283 *Eng Rev* **13**(1): 5–29.
- 1284 Voyant C, Notton G, Kalogirou S, Nivet M, Paoli C, Motte F, Fouilloy A. 2017.
1285 Machine learning methods for solar radiation forecasting: A review. *Renew*
1286 *Energy* **105**: 569–582. doi:10.1016/j.renene.2016.12.095.
- 1287 Waltz E, Waltz T. 2017. Principles and practice of image and spatial data fusion,
1288 in Liggins II M, Hall D, Llinas J, eds., *Handbook of multisensor data fusion*.
1289 CRC Press: pp. 109–134. doi:10.1201/9781420053098.
- 1290 Wang X, Xu J, Liu M, Wei Z, Bu W, Hong T. 2017. An ontology-based approach
1291 for marine geochemical data interoperation. *IEEE Access* **5**: 13364–13371.
1292 doi:10.1109/ACCESS.2017.2724641.
- 1293 WoRMS Editorial Board. 2020. World Register of Marine Species (WoRMS). Ac-
1294 cessed: 2020-01-03. <http://www.marinespecies.org>.
- 1295 Xiao G, Calvanese D, Kontchakov R, Lembo D, Poggi A, Rosati R, Za-
1296 kharyashev M. 2018. Ontology-Based Data Access: A Survey, in Lang J,
1297 ed., *Proceedings of the Twenty-Seventh International Joint Conference on*
1298 *Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*, pp.
1299 5511–5519. ijcai.org. doi:10.24963/ijcai.2018/777.
- 1300 Zalando Research. 2019. flair: A very simple framework for state-of-the-art NLP.
1301 Retrieved March. 21st, 2020. [https://github.com/flairNLP/](https://github.com/flairNLP/flair)
1302 *flair*.
- 1303 Zeng ML. 2008. Knowledge Organization Systems (KOS). *Knowl Organ* **35**(2-3):
1304 160–182. doi:10.5771/0943-7444-2008-2-3-160.
- 1305 Zhou L, Cheatham M, Krisnadhi A, Hitzler P. 2018. A Complex Alignment
1306 Benchmark: GeoLink Dataset, in Vrandečić D, Bontcheva K, Suárez-
1307 Figueroa MC, Presutti V, Celino I, Sabou M, Kaffee L-A, Simperl E, eds.,
1308 *The Semantic Web - ISWC 2018 - 17th International Semantic Web Con-*
1309 *ference, Monterey, CA, USA, October 8-12, 2018, Proceedings, Part II*, vol.
1310 11137, pp. 273–288. Springer. doi:10.1007/978-3-030-00668-6_17.